

DA  
1112  
1992  
④

寄	贈
新城靖氏	平成 年 月 日

並列分散応用プログラムを対象とした  
オペレーティング・システムの研究

1992年

新城 靖

# 目次

第1章	序論	1
1.1	並列／分散応用プログラムとオペレーティング・システム	1
1.2	資源割当てモジュール	3
1.3	競合と調和	5
1.4	主な研究成果	8
1.5	関連した研究	10
1.6	論文の構成	13
第2章	システムの構成とカーネルの機能	15
2.1	システムの概観	15
2.2	対象とするハードウェアの構成	16
2.2.1	結合モデルとサーバ・モデル	16
2.2.2	ネットワークと機密性	18
2.3	対象とする応用プログラム	18
2.3.1	並列／分散応用プログラムと逐次応用プログラム	18
2.3.2	利用者と応用プログラム	18
2.3.3	並列応用プログラムと分散応用プログラム	19
2.4	カーネル	19
2.5	オブジェクトとサーバ	21
2.6	オブジェクト識別子	21
2.6.1	オブジェクト識別子の構造	21
2.6.2	サイト識別子	22
2.6.3	位置情報を含むオブジェクト識別子に関する議論	24
2.7	遠隔手続き呼出し	25
2.7.1	位置一様性とサーバ型一様性	25
2.7.2	分散システムの多重視点	26
2.7.3	遠隔手続き呼出しを行うカーネル・コール	27
2.8	認証	28
2.8.1	利用者識別子	28
2.8.2	遠隔手続き呼出しとアクセス権の確認	28
2.9	名前サーバによる位置独立な名前付け	29

2. 1 0	カーネル・サーバ群	30
2. 1 0. 1	プロセス・サーバ	30
2. 1 0. 2	メモリ・サーバ	31
2. 1 0. 3	名前サーバ	31
2. 1 0. 4	ファイル・サーバ	31
2. 1 0. 5	タイム・サーバ	31
2. 1 1	関連した研究	31
2. 1 1. 1	カーネル化されたカーネル	31
2. 1 1. 2	マイクロカーネル	32
2. 1 2	まとめ	33
第3章	軽量プロセス	35
3. 1	軽量プロセスの概念	35
3. 2	軽量プロセス機能の実現目標	36
3. 3	代表的な軽量プロセス実現方式	37
3. 4	マイクロプロセスと仮想プロセッサ	38
3. 4. 1	他の実現方式との比較	40
3. 4. 2	プロセッサのスケジューリング	42
3. 5	マイクロプロセスを実現するライブラリの機能と構造	43
3. 5. 1	sleep/wakeup モデル	44
3. 5. 2	他の実現方式と比較	45
3. 5. 3	実行の軌跡の記録と解析	49
3. 6	仮想プロセッサを実現するカーネルの機能	52
3. 6. 1	カーネル・コール	52
3. 6. 2	カーネル・コールの利用例	52
3. 6. 3	仮想プロセッサの固有領域	55
3. 6. 4	仮想プロセッサの固有領域の利用	56
3. 6. 5	固有領域を利用しない方法との比較	57
3. 7	マイクロプロセスと仮想プロセッサの利用	58
3. 7. 1	共有メモリ型マルチプロセッサにおける並列処理	58
3. 7. 2	入出力の重ね合わせ	59
3. 7. 3	非同期通信の実現	61
3. 7. 4	動的負荷分散	62
3. 8	応用固有の同期・通信プリミティブの開発	63
3. 8. 1	層構造を持つライブラリの利用	63
3. 8. 2	コルーチンの性質の利用	64
3. 8. 3	応用固有の同期・通信プリミティブの開発例	65

3. 9	応用固有のスケジューラの開発	70
3. 9. 1	層構造を持つライブラリの利用	70
3. 9. 2	応用固有の軽量プロセス・スケジューラの開発例	70
3. 10	仮想プロセッサを実現するカーネルの構成法	72
3. 10. 1	仮想プロセッサの固有領域の内部実現	72
3. 10. 2	実プロセッサの固有領域	75
3. 10. 3	仮想プロセッサとプロセスのスケジューラ	76
3. 11	実現	79
3. 11. 1	仮想プロセッサを実現するカーネルの実現	79
3. 11. 2	マイクロプロセス・ライブラリの実現	81
3. 11. 3	UNIXにおける仮想プロセッサのエミュレーション	84
3. 12	性能	87
3. 12. 1	実験環境	87
3. 12. 2	比較対象	89
3. 12. 3	プロセスの生成・消滅	90
3. 12. 4	コンテキスト切替えを伴わないプロセス間の同期	93
3. 12. 5	コンテキスト切替えを伴うプロセス間の同期	95
3. 12. 6	仮想プロセッサの性能	98
3. 12. 7	応用固有のスケジューラの性能	98
3. 13	関連した研究	101
3. 13. 1	他の軽量プロセス実現方式との比較	101
3. 13. 2	他の仮想プロセッサ方式を用いる方式との比較	102
3. 13. 3	スピンロックとプロセッサの横取りの問題	103
3. 13. 4	共有 fork との比較	104
3. 13. 5	軽量プロセスの集合によるカーネルの構築	105
3. 14	まとめ	105
第4章	マッピング・コントローラ	107
4. 1	マッピング・コントローラ	107
4. 2	マッピング・コントローラに対する ReSC システムの方針	107
4. 3	応用固有のマッピング・コントローラの支援	108
4. 4	並列シェル	109
4. 4. 1	並列シェルの構文と機能	109
4. 4. 2	プロトタイプの実現	111
4. 4. 3	実験	112
4. 5	関連した研究	113
4. 5. 1	従来の分散型オペレーティング・システムとの比較	113

4. 5. 2	負荷均衡の方針の研究	114
4. 5. 3	負荷均衡の仕組みの研究	115
4. 6	まとめ	115
第5章	オブジェクトの堆積	117
5. 1	目標と背景	117
5. 2	基本概念	119
5. 2. 1	StdFS、ZFS、CFZ	122
5. 2. 2	サーバの堆積とオブジェクトの堆積	124
5. 2. 3	層化プログラミングとオブジェクトの堆積	125
5. 2. 4	UNIXのパイプラインとオブジェクトの堆積	125
5. 2. 5	抽象データ型、部分型とオブジェクトの堆積	126
5. 2. 6	C++によるオブジェクトの堆積の例	126
5. 2. 7	継承とオブジェクトの堆積	130
5. 2. 8	委譲とオブジェクトの堆積	131
5. 2. 9	メタ・オブジェクトとオブジェクトの堆積	131
5. 3	システムに対する要求	133
5. 4	オブジェクトの堆積による	
	分散型オペレーティング・システムの構築	135
5. 4. 1	望まれる性質	136
5. 4. 2	堆積を生成する手続き	137
5. 4. 3	ひな形オブジェクト・ベース	138
5. 4. 4	書込み時コピー	139
5. 5	種々の堆積可能オブジェクトとサーバ	139
5. 5. 1	フィルタ・サーバ	141
5. 5. 2	間接サーバ	142
5. 5. 3	グループ・サーバ	143
5. 6	分散型堆積可能オブジェクトとサーバ	144
5. 6. 1	複製ファイル・サーバ	144
5. 6. 2	間接オブジェクトによるオブジェクトの移動	145
5. 6. 3	分散型堆積可能オブジェクト間の競合	147
5. 7	多重視点	147
5. 7. 1	堆積の上下を参照することによる多重視点	148
5. 7. 2	多重視点とキャッシュ	149
5. 8	遠隔手続き呼出しスタブ生成器とインタフェース記述言語	150
5. 8. 1	インタフェース記述言語の構文と意味	150
5. 8. 2	SunRPCのインタフェース記述言語との比較	152

5. 9	堆積可能サーバの実現	154
5. 9. 1	典型的な堆積可能サーバの構成要素	154
5. 9. 2	キャッシングを利用した 融合ディレクトリ・サーバの実現	157
5. 9. 3	外部のオブジェクトへのキャッシュの保存	158
5. 9. 4	オブジェクトの堆積を実現する環境	158
5. 10	関連した研究	159
5. 10. 1	Plan 9	160
5. 10. 2	Interposition	160
5. 10. 3	多重継承	160
5. 10. 4	手続きの堆積	161
5. 11	まとめ	162

第6章	結論	165
-----	----	-----

謝辞	171
----	-----

参考文献	173
------	-----

付録	181
----	-----

# 図の目次

図 1 - 1	並列／分散ハードウェア環境	1
図 1 - 2	従来のシステムにおける 2 つの競合	6
図 1 - 3	ReSC システムの構成	8
図 2 - 1	対象とするハードウェアと応用プログラム	17
図 2 - 2	カーネル・サーバと外部サーバ	20
図 2 - 3	局所サイト型オブジェクト識別子による 位置独立のオブジェクトの識別	23
図 2 - 4	放送型、マルチキャスト型オブジェクト識別子による 位置独立のオブジェクトの識別	23
図 2 - 5	分散型システムの多重視点	27
図 2 - 5 (a)	分散型システムとして見る	27
図 2 - 5 (b)	集中型システムとして見る	27
図 3 - 1	プロセス、マイクロプロセス、仮想プロセッサの関係	38
図 3 - 2	軽量プロセス実現方式の比較	41
図 3 - 2 (a)	マイクロプロセス／仮想プロセッサ方式	41
図 3 - 2 (b)	カーネル制御方式	41
図 3 - 2 (c)	コルーチン方式	41
図 3 - 3	マイクロプロセス・ライブラリの構造	43
図 3 - 4	2 進セマフォの実現	46
図 3 - 4 (a)	マイクロプロセス・ライブラリによる	46
図 3 - 4 (b)	カーネル制御方式による	47
図 3 - 4 (c)	コルーチン方式による	48
図 3 - 5	実行の軌跡の記録より得られる統計情報	50
図 3 - 6	実行の軌跡の記録を視覚化したもの	51
図 3 - 7	vp_allocate() カーネル・コールによる仮想プロセッサの生成	53
図 3 - 8	vp_sleep() カーネル・コールによる実プロセッサの開放	54
図 3 - 9	vp_switch() カーネル・コールによる 他の仮想プロセッサへの制御の移動	54
図 3 - 1 0	2 つの仮想プロセッサの論理アドレス空間と物理メモリの対応	55

図 3 - 1 1	配列と仮想プロセッサ識別子による	57
図 3 - 1 2	ループ分割する前の逐次プログラム	59
図 3 - 1 3	マイクロプロセスによるループ分割	60
図 3 - 1 4	マイクロプロセスによる入出力の重ね合わせ	61
図 3 - 1 5	プロセスの移動	62
	図 3 - 1 5 (a) 重量プロセスの移動	62
	図 3 - 1 5 (b) 軽量プロセスの移動	62
図 3 - 1 6	共有メモリ型マルチプロセッサにおける マイクロプロセスを利用した S M A S H の実現	68
図 3 - 1 7	ネットワーク環境における マイクロプロセスを利用した S M A S H の実現	69
図 3 - 1 8	プロセス変換表と仮想プロセッサ変換表	73
図 3 - 1 9	2 段階のレディ・キュー	78
図 3 - 2 0	実験に用いた関係データベースへの問い合わせの構造	99
図 4 - 1	ストリームで結合されたコマンドのマッピング	111
	図 4 - 1 (a) 並列シェルによるもの	111
	図 4 - 1 (b) 簡単な負荷分散を行う マッピング・コントローラによるもの	111
図 5 - 1	オブジェクト、サーバ、クライアントの関係	119
図 5 - 2	サーバに管理されているオブジェクトと堆積オブジェクト	121
図 5 - 3	圧縮、暗号、標準ファイル・サーバ	123
図 5 - 4	サーバの堆積と層化プログラミング	124
	図 5 - 4 (a) オブジェクトの概念があるサーバの堆積	124
	図 5 - 4 (b) オブジェクトの概念がある層化プログラミング	124
図 5 - 5	C++ 言語によるオブジェクトの堆積の実現	127
	図 5 - 5 (a) 抽象データ型の定義	127
	図 5 - 5 (b) サーバのインタフェースの定義	127
	図 5 - 5 (c) サーバの実現	128
	図 5 - 5 (d) クライアント	128
図 5 - 6	main プログラムの実行結果	129
図 5 - 7	さまざまな堆積可能サーバ	141
図 5 - 8	複製ファイル・サーバ	145
図 5 - 9	間接サーバによるオブジェクトの移動	146
図 5 - 1 0	3 つの上位層のオブジェクトによる 3 つの視点の実現	148



図 5 - 1 1	堆積オブジェクトによる 3 つの視点の実現	148
図 5 - 1 2	同期型キャッシュ・サーバ	149
図 5 - 1 3	簡単なファイル・サーバのインタフェース記述	151
図 5 - 1 4	S u n R P C における 簡単なファイル・サーバのインタフェース記述	153
図 5 - 1 5	典型的な堆積可能サーバの構造	154
図 5 - 1 6	オブジェクト識別子とファイル構造体	156
図 5 - 1 7	外部のオブジェクトへキャッシュを保存する Z F S の実現	158

# 表の目次

表 1 - 1	2 つの調和と本論文の構成	13
表 3 - 1	プロセスの生成・消滅の実行時間の比較	91
表 3 - 1 (a)	1 回目	91
表 3 - 1 (b)	2 回目以降	92
表 3 - 2	コンテキスト切替えを伴わないプロセス間の同期の実行時間の比較	94
表 3 - 3	コンテキスト切替えを伴うプロセス間の同期の実行時間の比較	96
表 3 - 4	仮想プロセッサの生成・消滅の性能	98
表 3 - 5	仮想プロセッサのコンテキスト切替えの性能	98
表 3 - 6	データベースの並列処理の実行結果	100
表 4 - 1	実験に用いたスクリプトの性質	112
表 4 - 2	スクリプトの実行時間	113

# 第1章 序論

この論文では、並列応用プログラム、および、分散応用プログラムを対象としたオペレーティング・システムの設計と実現について述べる。この章では、はじめに、本研究の背景となった分散型オペレーティング・システム、共有メモリ型マルチプロセッサ用のオペレーティング・システム、および、並列／分散応用プログラムについて考察する。次に、本研究の目的を、調和という概念を用いて述べる。そして、本研究の主な研究成果、および、関連した研究について概略を述べる。最後に、本論文の構成について述べる。

## 1.1 並列／分散応用プログラムとオペレーティング・システム

共有メモリ型マルチプロセッサ、および、単一プロセッサが高速LAN (Local Area Network) によって結合された並列／分散ハードウェア環境が広く普及してきた (図1-1)。このようなハードウェア環境の利用は、エンジニアリングの分野に始まった。近年では、事務処理の分野においても、1台の大型計算機を共有する形態から、複数の個人用ワークステーションをネットワークで結合して利用する形態へと変化してきている。これは、大型計算機を複数人で共有するよりも、マイクロプロセッサを利用した個人用ワークステーションをネットワークで結合して利用の方が価格性能比がよくなってきたからである。共有メモリ型マルチプロセッサも、もはや珍しいも

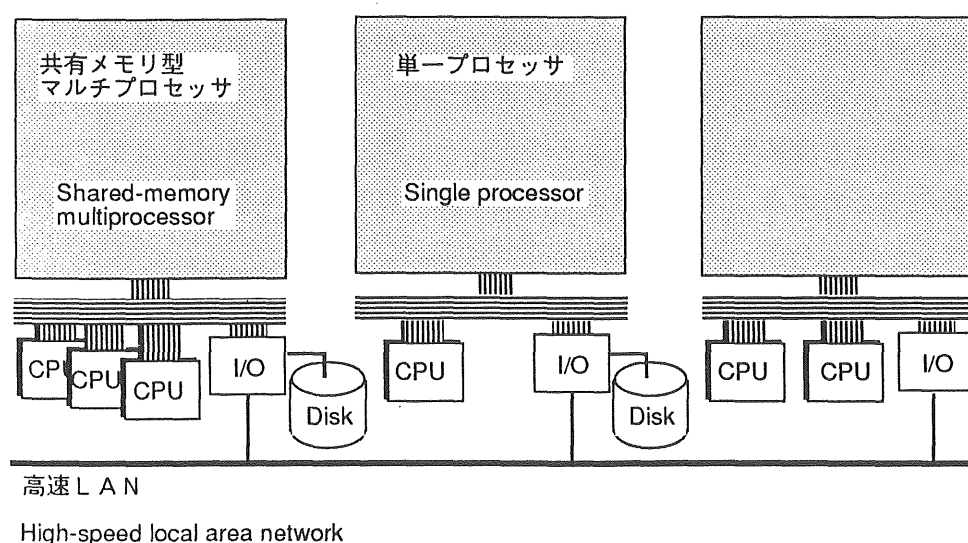


図1-1 並列／分散ハードウェア環境

Figure 1-1: A parallel/distributed hardware environment.

のではなく、しばしばワークステーションの上位モデルとして提供されている。

図1-1に示したような並列/分散ハードウェア環境の利用法を、次の2種類に分類することができる。

(1) 単一プログラミング・システム

(2) 多重プログラミング・システム

(1)では、そのハードウェア環境に含まれる計算機資源を1つの応用プログラムが専有し、その応用プログラムについて高い性能(performance)、信頼性

(reliability)、および、可用性(availability)を実現することを目標とする。

ここで計算機資源とは、CPU、メモリ、ファイル等である。今までに数多くの並列/分散応用プログラムが、単一プログラミング・システムとして開発されてきた。このようなシステムの例としては、SODA [37]、Amber [11]があげられる。

単一プログラミング・システムは、しばしば単一のプログラミング言語によって記述される[5]。このような並列/分散処理記述言語の例として、ARGUS

[49]、CONIC [51]があげられる。これらのシステムでは、プログラミング言語のレベルにおいて、明示的に並列処理や分散処理が記述される。このほかに、関数型言語、論理型言語、オブジェクト指向言語等によりプログラムを記述し、言語処理系により並列性を抽出する方法も研究されている。

単一プログラミング・システムでは、プロセッサやメモリなどの計算機資源は、全て言語処理系の管理下におかれる。言語処理系とは、コンパイラ、あるいは、実行時ルーチンのことである。コンパイラや実行時ルーチンは、並列/分散処理環境を最大限に活用するように、静的、あるいは、動的に資源割当ての最適化を行う。この場合、その応用の特性を利用することにより、資源割当ての最適化を容易に行うことができる。したがって、この方法は、1つの応用プログラムについてのみ高い性能と信頼性を実現する必要がある場合には、有効である。しかしながら、並列/分散ハードウェア環境を複数の応用プログラムにより共有する必要がある場合には、保護

(protection)の問題が生じるため、この方法を用いることはできない。さらに、他の応用プログラムと接続する必要がある場合にも、特別な工夫が必要となる。こうした要求に答えるためには、プログラムの保護や応用プログラム間のデータのやりとりを行うオペレーティング・システムが必要となる。

ネットワークにより結合された計算機群を、上記(2)の方法により利用するとは、分散型オペレーティング・システム(distributed operating system)を構築し、その上で複数の応用プログラムを走らせることである。代表的な分散型オペレーティング・システムの研究の例としては、Amoeba [54]、Locus [65]、Vシステム [12]があげられる。分散型オペレーティング・システムに関する重要な研究課題としては、ネットワーク上に分散したファイルやプリンタなど計算機資源に対するアクセスの実現、高速なプロセス間通信の実現、負荷分散、フォールト・トレラントなファイル・システムの構築があげられる。

従来の分散型オペレーティング・システムの主要な研究課題は、分散透明性 (network transparency) の実現であった。分散透明性とは、オペレーティング・システムが応用プログラムに対して仮想的な1台の集中型システムというモデルを提供し、応用プログラムがネットワークの構造を意識することなく、ネットワークに接続された計算機資源を利用することを可能にすることである。分散透明性が実現されたシステムでは、集中型システムにおいて開発された逐次応用プログラムを、全く変更することなく、並列/分散ハードウェア環境において利用することが可能となる。

L A Nにより結合された計算機の利用の研究が進む一方、L A N上の個々の計算機のマルチプロセッサ化が進んでいる。商用のマルチプロセッサとして最も普及しているものは、共有メモリ型マルチプロセッサである [6] [32] [100]。これは、複数のプロセッサ、1つの共有メモリ、および、入出力装置が、1つのバスで結合された構成をとる (図1-1)。プロセッサには、メモリへのアクセス速度を向上し、共有バスの負荷を抑えるために、キャッシュ・メモリが付加される。

初期の共有メモリ型マルチプロセッサは、主に従来の単一プロセッサのスループットを向上させるために利用された。従来の単一プロセッサ用の逐次応用プログラムを複数同時に走らせることで、単位時間当たりの処理可能な仕事の数を増やすことが可能となる。並列処理を行う場合、最初システムでは、従来のプロセスが用いられた。しかしながらこの方法には、プロセスの操作、すなわち、プロセスの生成・消滅やプロセス間の同期・通信のオーバーヘッドが大きいという問題があった。このため、しばしば並列処理の効果が打ち消されることがあった [103]。

この問題を解決するために、並列処理のための効率的なプロセスの研究が盛んに行われるようになった。このようなプロセスは、軽量プロセスと呼ばれる。軽量プロセス (lightweight process) とは、従来の、多重プログラミング・システムにおける資源割当てと保護の単位としてのプロセスとは異なる、並列処理の単位としてのプロセスである。軽量プロセスは、多重プログラミング・システムの言葉であり、単一プログラミング・システムにおける並列処理の単位としてのプロセスと対応する。単一プログラミング・システムでは、保護の単位としてのプロセスは、存在しない。

単一プログラミング・システムにおいて、並列処理の単位としてのプロセスは、その応用プログラムにより管理され、スケジュールされる。この場合、応用固有の情報を利用して、効率的なスケジューリングを行うことが可能となる。一方、多重プログラミング・システムにおいても、最近では、効率的な並列処理を行うためには、軽量プロセスのスケジューリングが重要であるとの認識が深まってきた。これを支援する方法が、重要な研究課題となっている。

## 1. 2 資源割当てモジュール

資源割当て (resource allocation) とは、C P U、メモリ、ファイル等の計算機

資源を、ある目的をもっていくつかの主体に分配することである。目的としては、次のようなことがあげられる。

- ・ 高速処理の実現
- ・ フォールト・トレランスの実現
- ・ システム全体のスループットの改善
- ・ 利用者間の公平の実現

1つのプログラム（1つのオペレーティング・システムを含む）の中で、資源割当ての最適化を行う部分を資源割当てモジュールと呼ぶことにする。資源割当てモジュールには、次のようなものがある。

（1）スケジューラ：

実行可能なプロセスの中から次に実行するプロセスを選択する。

（本論文では、スケジューリングという言葉とマッピングという言葉とを区別して用いる。スケジューリングとは、1プロセッサ内部においてプロセスにCPUを割り当てることを意味するものとする。マッピングとは、ネットワークに結合された複数のプロセッサに対して、プロセスやデータを割り当てることを意味するものとする。）

（2）マッピング・コントローラ（分散オブティマイザ）：

プロセスやデータ・ファイルをネットワーク上のどのプロセッサに割り当てるかを決定する。

（3）ページャ：

ページングに基づく仮想記憶システムにおいて、主記憶が不足した場合、どのページを2次記憶に追い出すかを決定する。

（4）ディスク・バッファ管理モジュール：

ディスク・ブロックのバッファリング（キャッシング）を行う場合、複数の入出力バッファの中から犠牲（victim）となるものを選択する。

（5）ディスク・ヘッドのスケジューラ：

複数の入出力要求の中で、ヘッドの位置をもとに、次に入出力を行うディスク・ブロックを決定する。

このような資源割当てモジュールは、旧来からオペレーティング・システムの重要な構成要素であった。現在でも、ハードウェア環境が大きく変化したにも関わらず、資源割当てモジュールがオペレーティング・システムの重要な構成要素であることには変りがない。ところが、最近の並列／分散応用プログラムの出現により状況が変化してきた。すなわち、資源割当てモジュールを内部に含むような応用プログラムが数多く開発されてきたのである。

逐次応用プログラムは、内部に処理単位を1つしか含んでいない。このため、資源割当てモジュールを持つ必要がない。これに対して、並列／分散応用プログラムは、内部に複数の処理単位を含んでいる。このため、資源割当てモジュールが必要となる。並列／分散応用プログラムの第1の目的は、高速処理や高信頼性の実現である。この

目的のためには、資源割当てが非常に重要となる。したがって、逐次応用プログラムと比較して、並列／分散応用プログラム特徴は、内部に資源割当てモジュールを含んでいる点にあるといえる。特に並列／分散応用プログラムの場合には、上記の資源割当てモジュールの中で、(1) スケジューラと(2) マッピング・コントローラが重要である。

### 1. 3 競合と調和

従来のオペレーティング・システムは、主に逐次応用プログラムを想定して機能が設計されていた。このため、並列／分散応用プログラムという、内部に資源割当てモジュールを含む応用プログラムに対する支援が不十分であった。従来の分散型オペレーティング・システムの主要な研究課題は、分散透明性の実現であった。分散透明性が実現されたシステムの上で並列／分散応用プログラムを走らせようとする、様々な問題が生じる。並列／分散応用プログラムは、並列処理や分散処理のために複数のプロセッサを必要とするが、従来の分散型オペレーティング・システムは、定義上、1つの仮想的な集中型システムというモデルしか提供しない。ネットワーク上の資源の位置も、並列／分散応用プログラムにとっては、資源割当てを行う際の重要な情報となるが、従来の分散型オペレーティング・システムは、この情報を覆い隠してしまう。

多重プログラミングの共有メモリ型マルチプロセッサにおいて、内部にスケジューラという資源割当てモジュールを含む並列応用プログラムを走らせる場合においても、従来の分散型オペレーティング・システムと似た問題に遭遇する。オペレーティング・システムのスケジューラは、複数の利用者や応用プログラム間の公平なCPU資源の利用を目的として、スケジューリングを行う。この方針が、しばしば1つの応用プログラムの内部の、並列処理の単位としてのプロセス(軽量プロセス)に対しても適用される。すなわち、1つの応用プログラム内部の軽量プロセスが公平に実行される。これに対して、並列応用プログラムのレベルでは、公平なスケジューリングではなく、むしろ特定の軽量プロセスに多くのCPU資源を割り当て、その応用プログラムの全体の処理時間を短縮することを試みる。

分散型オペレーティング・システムでは、システムと並列／分散応用プログラム間でマッピングの方針に食い違いが生じる。同様に、共有メモリ型マルチプロセッサのオペレーティング・システムにおいても、システムと並列応用プログラム間でスケジューリングの方針に食い違いが生じる。このような食い違いを、オペレーティング・システムと並列／分散応用プログラム間の競合(conflict)と呼ぶことにする。これを、図1-2のConflict(1)に示す。

図1-2において、逐次応用プログラム(sequential application)は、オペレーティング・システムから分散透明性やマッピングという機能が提供されることを望ん

でいる。これに対して、並列／分散アプリケーション (parallel/distribute application) は、そのようなシステムの機能が不用であるばかりでなく、自分自身の資源割当てを妨害するものであると考える。このように、オペレーティング・システムが提供するある機能が、逐次アプリケーションにとっては必要であり、並列／分散アプリケーションにとっては、不用であることがある。このように、逐次アプリケーションと並列／分散アプリケーションの間で、システムに対する要求が対立することを、逐次アプリケーションと並列／分散アプリケーションの間の競合とよぶ。これを図1-2の Conflict(2) に示す。

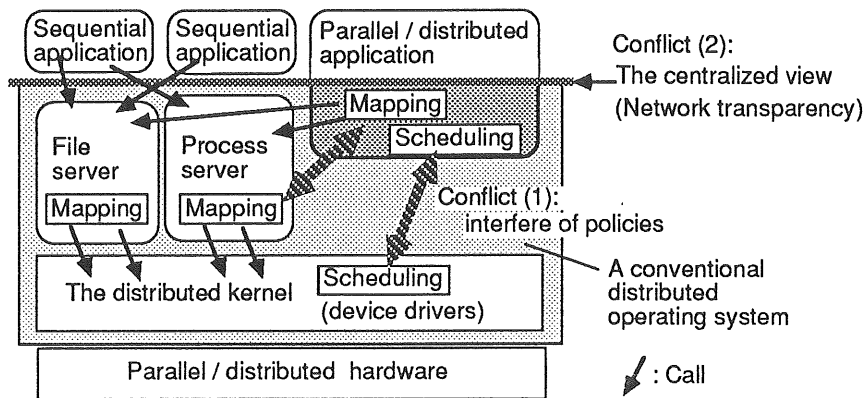


図1-2 従来のシステムにおける2つの競合

Figure 1-2: Conflicts in a conventional operating system.

本研究では、逐次アプリケーションに加えて、並列／分散アプリケーションという、内部に資源割当てモジュールを含むプログラムを対象とした分散型オペレーティング・システムにの設計と実現を行った [76～88]。このシステムを ReSC と名付けた。ReSC が対象とするハードウェア環境は、共有メモリ型マルチプロセッサ、および、単一プロセッサが高速ネットワークにより結合されたものである (図1-1)。ReSC が対象とするアプリケーションは、並列／分散アプリケーション、および、逐次アプリケーションである。ここで、並列／分散アプリケーションとは、内部に複数の処理単位と資源割当てモジュールを含むものとする。逐次アプリケーションは、内部に1つの処理単位しかなく、資源割当てモジュールを含んでいないものとする。

ここで導入した並列／分散アプリケーションと逐次アプリケーションの分類は、本システムの機能を説明するためのものである。したがって、一般的な意味における並列／分散アプリケーションが本システムの逐次アプリケーション用の機能を利用する場合も、その逆の場合も在り得る。また、システムの支援の観点では、並列アプリケーションと分散アプリケーションの差異は、小さい。なお、プログラマにより逐次的に記述されたプログラムが言語処理系により複数の処理単位に分割されたものは、本システムでは、並列／分散アプリケーションに分類される。



本研究の目的は、次の2種類の調和 (harmonization) を実現することに集約される [76] [79]。

(1) オペレーティング・システムと並列/分散応用プログラム間の調和:

システムは、並列/分散応用プログラムと競合することなく、そのような応用プログラムに対して、高速処理や高信頼性を実現するためにスケジューリングやマッピングを制御する機能を提供する。

(2) 逐次応用プログラムと並列/分散応用プログラム間の調和:

1つの並列/分散ハードウェア環境の中で、並列/分散応用プログラムと逐次応用プログラムを同時に利用することを可能にする。逐次応用プログラムが、システムを集中型システムとして見ることを可能にする。これにより、集中型システムで開発された逐次応用プログラムを一切変更することなく、ネットワークで結合された計算機群で実行することが可能となる。同時に、並列/分散応用プログラムがシステムを並列/分散システムとして見ることを可能にする。これにより、並列/分散応用プログラムは、並列/分散ハードウェア環境をほとんどそのままの形で利用することができる。

上記(2)の目標の、逐次応用プログラムに関する部分は、従来の分散型オペレーティング・システムの研究の目標と一致する。ただし、応用プログラムの実行効率に関しては、内部に資源割当てモジュールを含む並列/分散応用プログラムを優先することにする。この結果、逐次応用プログラムは、従来の分散型オペレーティング・システムの方が効率的に動作することが予想される。しかしながら、システムの資源を大量に消費する並列/分散応用プログラムの実行効率を改善することが、結果としてシステム全体の処理能力の改善につながると思われる。

上記の目標を実現するために、分散型オペレーティング・システムを、オペレーティング・システム・カーネル上で動作する並列/分散応用プログラムの集合として設計した。これを、図1-3に示す。システムは、カーネル(kernel)、外部サーバ群(external servers)、並列シェル(parallel shell)から構成される。カーネルは、プロセス間通信等の基本的なサービスと、並列/分散応用プログラムのための効率的なサービスを提供する。外部サーバ群は、ファイルやディレクトリのキャッシング、複製、および、移動といった高度な機能を提供する。並列シェルは、ネットワーク上の負荷が小さいプロセッサを探し、コマンド群を並列に実行する。

システムの構成要素である外部サーバ群と並列シェルは、並列/分散応用プログラムと同じ立場でカーネルが管理している計算機資源を利用する。これにより、上記の調和(1)が実現される。すなわち、並列/分散応用プログラムは、システムの要素である外部サーバ群や並列シェルに含まれている資源割当てモジュールと競合することなく、カーネルが提供するシステムの計算機資源を利用することができる。

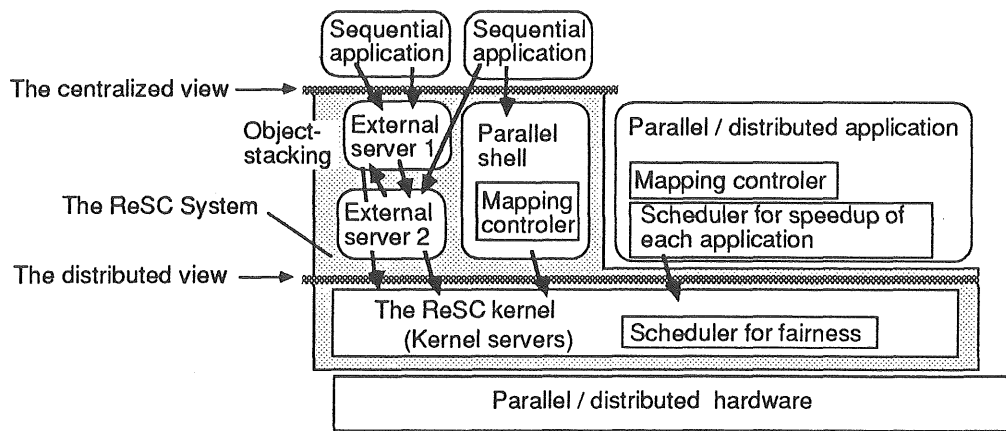


図 1-3 ReSCシステムの構成

Figure 1-3: The organization of the ReSC system.

上記の調和(2)は、並列シェルと外部サーバ群により実現される。逐次応用プログラムは、並列シェルと外部サーバ群によって提供される、従来の分散型オペレーティング・システムと同じ環境を利用することができる。同時に、並列/分散応用プログラムは、これらの機能を利用することなく、カーネルを直接利用することで、並列/分散ハードウェア環境を(多重プログラミングの範囲内で)ほとんどそのままの形で利用することができる。

本研究は、RISC (reduced instruction set computer) アーキテクチャの影響を受けている。本システムの ReSC という名前は、a distributed operating system with Reduced System Calls の略で、本来は、オペレーティング・システム全体ではなくカーネルの名前であった。(現在では、「ReSCシステム」として、オペレーティング・システム全体を指すためにも用いる。) RISCでは、機械語命令のレベルが内部のALUのレベルに設定されており、コンパイラがパイプラインや遅延スロットを考慮した最適なコードを生成することを期待している。これに対して、ReSCでは、システム・コールのレベルが下位層の並列/分散ハードウェアのレベルに設定されており、並列/分散応用プログラムがマッピングやスケジューリングの最適化を行うことを期待している。

#### 1.4 主な研究成果

本研究では、並列/分散応用プログラム、および、逐次応用プログラムを対象としたオペレーティング・システムに関連して、幅広い範囲で研究を行った。その結果得られた成果を以下にまとめる。

### (1) システムの全体構成の提案

本研究では、1.3節で述べた2つの調和を実現するという研究の目的に適したオペレーティング・システムの全体構成について提案を行った[76][79]。オペレーティング・システムは、全体では様々な機能を提供するモジュールの巨大な集合である。したがって、それらのモジュールを整理する方法、ならびに、応用プログラムに対するインタフェースは、重要な研究課題となっている。具体的には、次のような構成要素に対する機能分割の方法とインタフェースが課題となっている。

- ・カーネル
- ・カーネル外のサーバ
- ・ライブラリ
- ・シェル
- ・コマンド群
- ・公のファイル

本研究では、これらの構成要素に対する機能の分割の方法を提案した。

### (2) カーネル機能の提案

システムの構成要素の中で特に重要となるカーネル機能の提案を行った[82][87]。カーネルは、他の構成要素とは異なり、オペレーティング・システムの中で唯一実行時に利用者定義のものに取り替えられない部分である。ゆえに、どのような機能をカーネルに持たせるかが重要な研究課題となっている。本研究では、カーネルを、全ての応用プログラムに対して機能を提供するモジュールであり、多重プログラミング・システムとしての資源割当てを行う部分として位置付けた。一方、並列／分散応用プログラムにとっては、カーネルがその動作環境全体となることを考慮して機能を設計した。さらに、応用プログラムの性質に応じて、システムが集中型システムとしても分散型システムとしても利用可能になるように、システム・レベルのオブジェクト識別子とプロセス間通信機能を設計した。

### (3) 軽量プロセス実現方式の提案

応用固有の同期・通信プリミティブとスケジューラの開発を支援する軽量プロセス実現方式の提案を行った[81][82][84][87][88]。具体的には、利用者レベルの軽量プロセスであるマイクロプロセスとカーネル・レベルの実体である仮想プロセッサにより軽量プロセス機能を提供する方法を提案した。マイクロプロセスは、利用者レベルで効率的に実現される。さらに、マイクロプロセスは、層構造を持つライブラリにより提供されるため、応用固有の同期・通信プリミティブやスケジューラが容易に実現される。

### (4) 仮想プロセッサを提供するカーネルの構成法の提案

仮想プロセッサを提供するオペレーティング・システムのカーネルを、カーネル上

で動作する並列応用プログラムと同じ手法を用いて構成する方法を提案した [77]。具体的には、利用者に対する抽象である仮想プロセッサをカーネル内の軽量プロセスとして実現し、その軽量プロセスを実プロセッサで実行する方法を提案した。その結果、カーネルが応用プログラムと相似形になり、さまざまな利点が生じる。また、この構成法は、他の外部仕様を持つ仮想プロセッサやカーネル・レベルの軽量プロセスを実現する際にも利用可能である。

#### (5) オブジェクトの堆積モデルの提案

object-basedシステムにおける複数のオブジェクトの機能を統合して利用するためのモデルとして、オブジェクトの堆積 (object-stacking) を提案した [78]

[80] [83] [85]。そして、このモデルを利用するためのシステムの条件を明らかにした。さらに、分散型オペレーティング・システムにおける有用なサーバを提示した。

### 1.5 関連した研究

この節では、本研究の動機となった主な研究について述べる。具体的な他の研究との比較については、個々の章で述べる。

#### (1) 従来の分散型オペレーティング・システム

文献 [98] において、Tanenbaum は、次のように分散型オペレーティング・システムを定義している。

「利用者には、通常の集中型のオペレーティング・システムのように見えるが、複数の独立したCPU (ネットワーク上のプロセッサ) の上で実行されるオペレーティング・システムである。・・・利用者は、システムを異なった計算機としてではなく、“仮想的な1台のプロセッサ”として認識する。・・・」

その他の多くの文献においても、この文献が引用され、さらに、ネットワーク・オペレーティング・システムと対比することで、類似の定義を行っている [22]

[43]。

1.3節で述べたように、この定義に従った分散型オペレーティング・システム上では、並列/分散応用プログラムが動作しない。すなわち、“仮想的な1台のプロセッサ”では、並列処理を行うことも、分散処理を行うことも不可能である。実際のシステムでは、現実的な選択として、特殊な方法によりネットワークの位置の情報を提供している。

本研究では、この分散型オペレーティング・システムの定義、すなわち、分散型オペレーティング・システムの研究の目標を新たに設定し直した。研究の目標は、

1. 3節で述べた2つの調和を実現することである。本研究における分散型オペレーティング・システムの定義を、以下に示す。

分散型オペレーティング・システムとは、ネットワークにより結合された計算機群を統一的に管理するオペレーティング・システムである。個々の計算機で独立したオペレーティング・システムが動作するのではなく、全体で1つのオペレーティング・システムが動作する。1つの分散型オペレーティング・システム上では、共通のオブジェクト識別子や利用者識別子が利用可能である。そして、ネットワークに接続された全ての計算機において、統一的なアクセス権のチェックが行われる。また、ネットワーク上の異なるプロセッサに配置されたプロセス間の通信も、同一のプロセッサに配置されたプロセス間の通信と全く同一の方法で実行することができる。

ネットワーク通信機能を持ちながら、個々の計算機で独立したオペレーティング・システムが動作するものは、ネットワーク・オペレーティング・システム (network operating system) とよばれている [98] [22] [43]。この定義は、本研究においても有効である。

## (2) 並列／分散応用プログラムを対象とした単一プログラミング・システム

並列／分散応用プログラムを対象とした代表的な単一プログラミング・システムとしては、SODA [37] と Amber [11] があげられる。SODA (a Simplified Operating system for Distributed Applications) は、本研究と同様に、RISCアーキテクチャの影響を受けている。SODAでは、多重プログラミング・システムの複雑さを排除することにより、分散応用プログラムに適したオペレーティング・システムを構築することを目指している。本研究では、従来の分散型オペレーティング・システムと同様に、多重プログラミングのシステムを構築することを目標とする。

Amberは、ワシントン大学で開発された、高い性能を得ることを目標とする応用プログラムのためのプログラミング・システムである [11]。Amberは、ネットワークで結合された計算機群において、移動可能オブジェクト (mobile object) とスレッドからなる、object-basedプログラミング・モデルを提供する。Amberのプログラマは、オブジェクトの移動 (object migration) を使ってネットワーク上のデータとプロセスの位置を制御する。

Amberでは、応用プログラムがネットワークの構成を知り、その応用プログラムの最大限の高速化のためにそれを利用する。すなわち、ネットワーク上のデータの配置は、実行時ルーチン (Amberシステム) ではなく、プログラマにより制御されるべきであると主張している。ReSCシステムも、Amberと全く同じ立場である。相違点としては、Amberは、単一プログラミング・システムであるのに対して、ReSCは、多重プログラミング・システムを目指していることがあげられる。

A m b e r では、C ++を制限したオブジェクト指向言語を提供しているのに対して、R e S C では、特にプログラミング言語を規定していない。

### (3) 軽量プロセスに関する研究

軽量プロセスが普及する以前、共有メモリ型マルチプロセッサにおける並列処理では、従来の重たいプロセスが利用されていた。この段階においても、システムの負荷に適応 (adapt) させるために、実行時にプロセスの数を変更することか試みられている [103]。一方、単一プロセッサでは、コルーチンによる軽量プロセスの実現が数多く行われた [93] [106]。コルーチンの場合、並列処理を行うことができないという大きな問題点がある。

共有メモリ型マルチプロセッサにおける軽量プロセスの研究で最初に大きな注目を集めたシステムは、カーネギー・メロン大学で開発された M a c h システムである [2] [67]。M a c h では、従来のプロセスに相当するタスクに加えて、スレッドとよばれる軽量プロセスを提供している。

M a c h では、カーネルにおいて軽量プロセスの実現を行っている。これに対して、利用者レベルにおいて軽量プロセスを実現し、それをカーネル・レベルの仮想プロセッサにより実行する方式が、本研究 [88] [81] をはじめ、数多く提案された [3] [23] [30] [53]。これらの方式の特徴は、カーネル・レベルの軽量プロセスと比較してよりオーバーヘッドが小さいこと、および、利用者レベルの情報が得やすいこと等があげられる。これらのシステムとの詳細な比較は、3.13節において行う。

### (4) オブジェクト指向システムとオブジェクト指向プログラミング言語

最近の分散型オペレーティング・システムの多くは、object-basedシステムとして構成されている [98] [22] [43]。object-basedシステムとは、継承 (inheritance)、あるいは、委譲 (delegation) の機能が支援されている必要がないオブジェクト指向システムである。オペレーティング・システムのレベルで継承、あるいは、委譲を支援している例は、非常に稀である [43]。したがって、オブジェクト指向オペレーティング・システムとよばれているシステムの多くは、object-basedシステムに分類される [1] [15] [109]。

本研究では、object-basedシステムにおいて、複数のサーバを統合して利用するためのモデル化の手法として、オブジェクトの堆積を提案する。UNIXのパイプは、単純なコマンドを結合し、複雑なコマンドを作るために使われる [4] [46]。これに対してオブジェクトの堆積は、単純なオブジェクトを結合し、複雑なオブジェクトを作るために使われる。層化プログラミングは、再利用可能性が高いモジュールを構築するために有力な方法の1つである [16]。同様に、オブジェクトの堆積も、再利用可能性が高いモジュールをサーバとして開発するための有効な方法である。

## 1.6 論文の構成

1.3節で述べた2つの調和と本論文の構成の関係を、表1-1にまとめる。この表においては、逐次応用プログラムと並列/分散応用プログラムの間の調和として、逐次応用プログラムに対して提供する機能を示している。正確には、この調和は、逐次応用プログラムがここであげた機能を利用し、並列/分散応用プログラムがカーネルを直接利用することで実現される。

表1-1 2つの調和と本論文の構成

章 \ 調和	システムと並列/分散応用プログラムの間の調和	逐次応用プログラムと並列/分散応用プログラムの間の調和 (逐次応用プログラムに対する機能)
2章 カーネル	<ul style="list-style-type: none"> <li>・カーネル・サーバ群</li> <li>・位置情報を含むオブジェクト識別子</li> </ul>	<ul style="list-style-type: none"> <li>・名前サーバ</li> <li>・位置独立のRPC発行</li> </ul>
3章 軽量プロセス	<ul style="list-style-type: none"> <li>・利用者レベルの軽量プロセス</li> </ul>	<ul style="list-style-type: none"> <li>・仮想プロセッサの自動生成</li> </ul>
4章 マッピング コントローラ	<ul style="list-style-type: none"> <li>・応用固有のマッピングコントローラの支援</li> </ul>	<ul style="list-style-type: none"> <li>・並列シェル</li> </ul>
5章 オブジェクト の堆積		<ul style="list-style-type: none"> <li>・外部サーバ群による高度な機能</li> </ul>

本論文は、以下のように構成されている。

第2章では、システムの構成とカーネルの機能について述べる。まず、本研究の対象となるハードウェア、および、応用プログラムについて深く考察する。次にカーネルの機能とカーネル・レベルのオブジェクト識別子、および、カーネル内部のサーバについて述べる。

第3章では、軽量プロセスについて述べる。まず、軽量プロセスの概念を明確に規定し、本システムにおける軽量プロセス実現方式であるマイクロプロセス/仮想プロ

セッサ方式について述べる。次に、これを利用した応用固有の同期・通信プリミティブとスケジューラの開発について述べる。続いて、仮想プロセッサを提供するカーネルの構成法について述べる。そして、実現した軽量プロセス機能の性能を評価し、他のシステムとの比較を行う。

第4章では、マッピング・コントローラについて述べる。始めに、マッピング・コントローラ概念を明確にする。続いて、ReSCシステムにおけるマッピング・コントローラの支援について述べる。最後に、逐次応用プログラムのためのマッピング・コントローラである並列シェルの機能とそのプロトタイプ性能について述べる。

第5章では、オブジェクトの堆積について述べる。まず、オブジェクトの堆積のモデルとそれを利用するためのシステムへの要求を整理する。次に、分散型オペレーティング・システムにおけるオブジェクトの堆積の利用について述べる。

第6章では、本論文のまとめと今後の研究の展開について述べる。



## 第2章 システムの構成とカーネルの機能

この章では、R e S C システムの構成とカーネルの機能について述べる。R e S C は、多重プログラミングの object-based 分散型オペレーティング・システムである。R e S C は、並列／分散応用プログラム、および、逐次応用プログラムの両方を支援する。システムの機能は、カーネル、カーネル外のサーバ、ライブラリ、シェルにより提供される。R e S C の特徴は、カーネルの機能選択、および、システム・レベルのオブジェクト識別子に現れる。

本章では、まず本研究の対象としたハードウェア、および、応用プログラムについて述べる。次にカーネルの機能とシステム・レベルのオブジェクト識別子、および、カーネル内部のサーバについて述べる。最後に、関連した研究との比較を行う。

本システムが提供する機能の中で、軽量プロセスに関する機能については、第3章において述べる。マッピング・コントローラを支援する機能、および、シェルにより提供する機能については、第4章において述べる。第5章で述べるオブジェクトの堆積の実現では、本章で述べるカーネルの機能が重要な役割を果たす。オブジェクトの堆積で必要とされる機能については、第5章において詳しく述べる。

### 2.1 システムの概観

本R e S C システムは、共有メモリ型マルチプロセッサ、および、単一プロセッサがLANにより結合されたハードウェアを管理するオペレーティング・システムである。従来の分散型オペレーティング・システムと同様に、本システムは、マルチプログラミング、かつ、マルチユーザのシステムである。本システムでは、複数の逐次応用プログラム、および、複数の並列／分散応用プログラムが同時に動作する。そして、複数の利用者が同時にログインし、システムの機能を利用する。

システムは、カーネル、外部サーバ、ライブラリ、並列シェルから構成される。カーネルは、特権モードで動作し、下位層の並列／分散ハードウェアを管理する。そして、応用プログラムに対してハードウェアの能力を提供する。

並列応用プログラムに対しては、カーネルは、下位層の並列／分散ハードウェアをほとんどそのままの形で提供する。このような応用プログラムからは、並列／分散ハードウェアを単独で使っているように見えるかもしれない。しかしながら、本システムは、多重プログラミング・システムであり、プロセッサ、メモリ、ネットワーク等の資源は、全ての応用プログラムにより共有されている。本システムのカーネルは、これらの資源を管理し、公平な資源の割当てを行う。

逐次応用プログラムに対しては、従来の分散型オペレーティング・システムと同様に、分散透明な資源参照を可能にする。そのために、R e S Cシステムは、ネットワーク上の位置に依らない遠隔手続き呼出しの発行、および、名前サーバによる位置に依らないオブジェクトの名前付けの機能を提供する。

本システムでは、様々なサービスを提供するプロセスをサーバとよぶ。カーネルは、クライアント (client) とサーバの間の通信機能を提供する。サーバは、カーネル内のサーバとカーネル外のサーバに分類される。前者を、カーネル・サーバ (kernel server)、または、内部サーバ (internal server) とよぶ。後者を、外部サーバ (external server) とよぶ。カーネル・サーバは、プロセス管理、メモリ管理、ファイル入出力、ネットワーク通信等の基本的な機能を提供する。カーネル・サーバは、プロセスではないが、クライアント・プロセスからは、遠隔手続き呼出しという、プロセスと同じインタフェースによりアクセスすることが可能になっている。より高度な機能は、外部サーバにより提供される。並列／分散応用プログラムは、主にカーネル・サーバを利用し、逐次応用プログラムは、主に外部サーバを利用する。外部サーバも、1つのクライアント・プロセスとして、カーネル・サーバや別の外部サーバの機能を利用する。

## 2. 2 対象とするハードウェアの構成

図2-1に、R e S Cが対象とするハードウェアを示す。これは、L A Nで結合された複数の計算機から構成されている。各計算機を、サイト (site) とよぶことにする。各サイトは、共有メモリ型マルチプロセッサ、または、単一プロセッサである。各サイトの構成、すなわち、プロセッサの数、メモリの容量、ディスクの有無は、それぞれ異なってもよい。

### 2. 2. 1 結合モデルとサーバ・モデル

文献 [64] において、Popek らは、システムの構成法を次の2種類に分類している。

#### (1) 結合モデル (integrated model)

ほぼ同質の計算機を対等に結合する。

#### (2) サーバ・モデル (server model)

ファイル・サーバ、プリント・サーバ、広域ネットワーク・サーバなど、専用機能を持つサーバを用いる。

この分類法に従うと、本システムが対象とするハードウェア構成は、(1) に分類される。A m o e b a [54] やVシステム [12] のハードウェア構成は、(2) である。

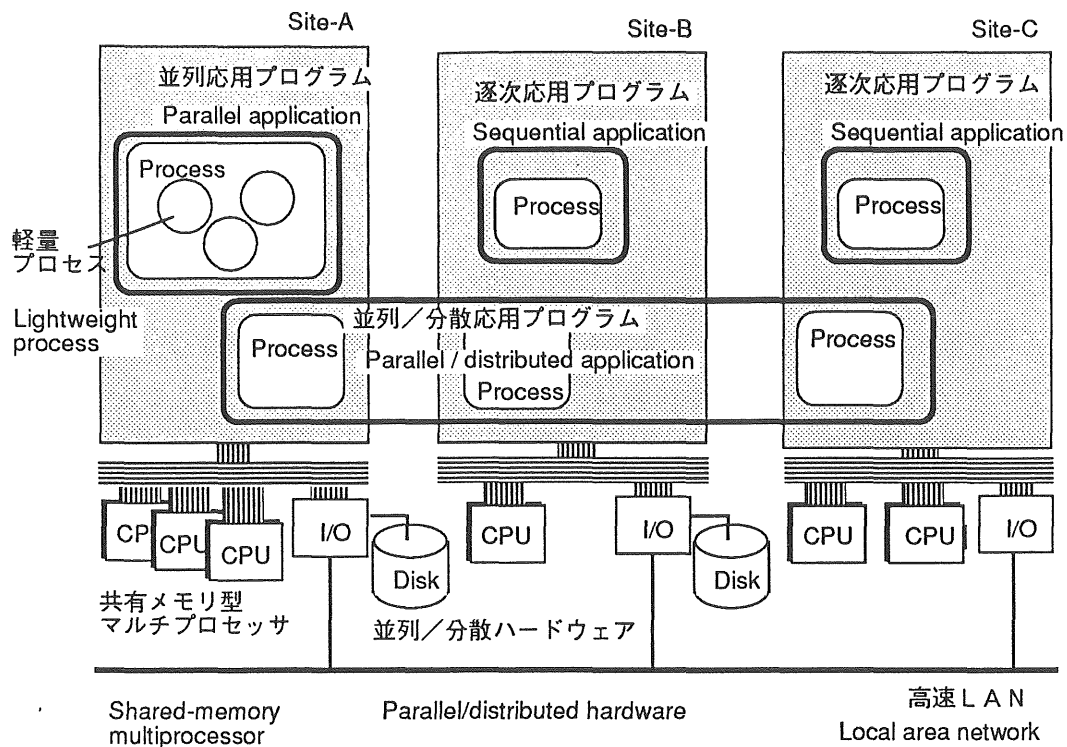


図 2-1 対象とするハードウェアと応用プログラム

Figure 2-1: The hardware environment and applications of ReSC.

このようなハードウェア構成を仮定した理由を、以下に示す。

(1) 入出力の高速化を図ることができる。近年、プロセッサの高速化とメモリの大容量化は、目をみはるものがある。これに対して、ディスク装置の速度は、それほど改善されていない。図 2-1 のように、複数の計算機にディスクを接続することにより、複数の入出力を並列に実行することが可能となり、入出力の高速化を図ることができる。また、ディスク・アレイ [62] と同様に、複数のサイトのディスクにデータを重複して保存することにより、ディスクの故障やサイトの故障に対応することが可能となる。

(2) 計算機の入替えが容易である。古い計算機を取り外したり、新しい計算機を組み入れることが容易である。大容量のディスクを供えた集中化されたファイル・サーバ専用計算機は、導入した瞬間には、すばらしい構成に思える。しかしながら、それを導入した翌年には、個人用のワークステーションがそのファイル・サーバよりも高速のプロセッサを供えることになる。さらに、そのファイル・サーバよりも大容量のディスクを備える可能性もある。

本システムでは、共有メモリ型マルチプロセッサの各プロセッサを、互に等価なも

のとして扱う。したがって、応用プログラムは、特定のプロセッサを指定してプロセスを生成することはできない。これは、キャッシュ・メモリの効果を見れば、各プロセッサから共有メモリへのアクセス時間が等しくなるからである。ただし、各応用プログラムに対して割り当てられた（複数の）プロセッサは、その応用プログラムの内部で自由に利用される。

## 2. 2. 2 ネットワークと機密性

本研究が対象とするネットワークは、LAN (Local Area Network) とする。WAN (Wide Area Network) との接続は、ゲートウエーを介して行うこととする。したがって、本研究では、ネットワークの機密性が保たれているものと仮定する。

## 2. 3 対象とする応用プログラム

### 2. 3. 1 並列／分散応用プログラムと逐次応用プログラム

ReSCでは、対象とする応用プログラムを次の2つに分類する（図2-1）。

#### (1) 並列／分散応用プログラム (parallel/distributed application programs)

内部に複数の処理単位と資源割当てモジュールを含んでいる。

#### (2) 逐次応用プログラム (sequential programs)

内部に1つの処理単位しか含んでいない。

この分類方法は、システムの支援の違いを議論するときに有用である。すなわち、本システムでは、前者に対しては、プログラムの実行効率を重視した支援を行い、後者に対しては、プログラムの記述性を重視した支援を行う。ただし、この分類は、完全に排他的ではない。ある並列／分散応用プログラムが、効率に大きな影響を与えない部分では、逐次応用プログラムに対する機能を利用することもある。

### 2. 3. 2 利用者と応用プログラム

ReSCでは、利用者 (user) と応用プログラム (application program) をはっきり区別して考える。利用者は、応用プログラムを利用して仕事をする人間である。応用プログラムは、利用者とは異なる人間 (プログラマ) により記述されるものと仮定する。この論文で述べる応用プログラムに対する支援とは、応用プログラムのプログラマに対して提供するシステムのモデル、並びに、そのプログラマが記述した実行時ルーチンに対して提供するシステムの機能を意味する。

### 2. 3. 3 並列応用プログラムと分散応用プログラム

システムによる支援の視点からは、並列応用プログラムと分散応用プログラムの差異は、小さい。したがって、本論文では、並列応用プログラムと分散応用プログラムをまとめて扱うことが多い。あえてこれらの2つを分類すると、次のようになる。

#### (1) 並列応用プログラム：

高性能 (high performance) を得ることを主目的とする応用プログラム。そのためには、共有メモリ型マルチプロセッサを利用したり、あるいは、ネットワークにより結合された計算機群を利用する。

#### (2) 分散応用プログラム：

高い信頼性 (reliability) や可用性 (availability) を得ることを主目的とする応用プログラム。そのためには、ネットワークに結合された複数の計算機を利用する。1台の共有メモリ型マルチプロセッサだけを利用することはない。

## 2. 4 カーネル

ReSCシステムは、カーネル、外部サーバ、ライブラリ、並列シェルから構成される (図2-2)。この中でカーネルは、最も重要な構成要素である。それは、全ての計算機上にコピーが存在し、全ての応用プログラムから利用されるからである。カーネルは、特権モードで動作する。外部サーバや並列シェル等の他の構成要素とは異なり、カーネルを利用者定義のものに取り替えることはできない。

カーネルは、以下のような機能を提供する。これ以外の機能は、外部サーバ、並列シェル、または、ライブラリの形で提供される。

#### (1) 共有資源の管理と公平な配分。

(メモリ管理、CPUのスケジューリング)

#### (2) プロセス間通信と保護。

#### (3) 並列/分散応用プログラムのための、単純で効率的な形の下位層のハードウェアが持つ能力。

(ファイル・サービス、プロセス・サービス、ネットワーク・サービスなど)

#### (4) 標準的なインタフェース

上記(1)と(2)は、多重プログラミング・システムとしての機能である。上記(3)に本システムの特徴がある。ReSCシステムは、他の分散型オペレーティング・システムと同様に、「カーネルを可能な限り小さくする」という方針に基づいて設計されている。相違点は、ReSCでは、カーネルがデバイス・サーバではなく、完全なサーバを含んでいる点にある。たとえば、Vシステムでは、カーネル中のデバ

イス・サーバは、それぞれの物理的なディスク装置について生のブロック型デバイス (a raw block device) としてのアクセスを提供する。そして、カーネル外のファイル・サーバは、そのデバイス・サーバを利用して、ファイルという抽象を実現する。R e S Cでは、カーネル中に完全なファイル・サーバが存在する。この設計により、カーネルが大きくなってしまいが、効率は、良くなる。なぜならば、保護の壁を横切るオーバーヘッドが減るからである。Vシステムとの違いに関しては、2. 1 1. 1項においてさらに詳しく述べる。

上記(4)は、オブジェクトの堆積を用いてサーバを統合する時に必要である。オブジェクトの堆積の詳細については、第5章において述べる。

カーネル内のサーバを、カーネル・サーバ (kernel server)、または、内部サーバ (internal server) とよぶ (図2-2)。一方、カーネル外にあるサーバを外部サーバ (external server) よぶ。外部サーバは、プロセスまたは、プロセスの集合である。カーネル・サーバについては、2. 1 0節において述べる。外部サーバについては、5章において述べる。

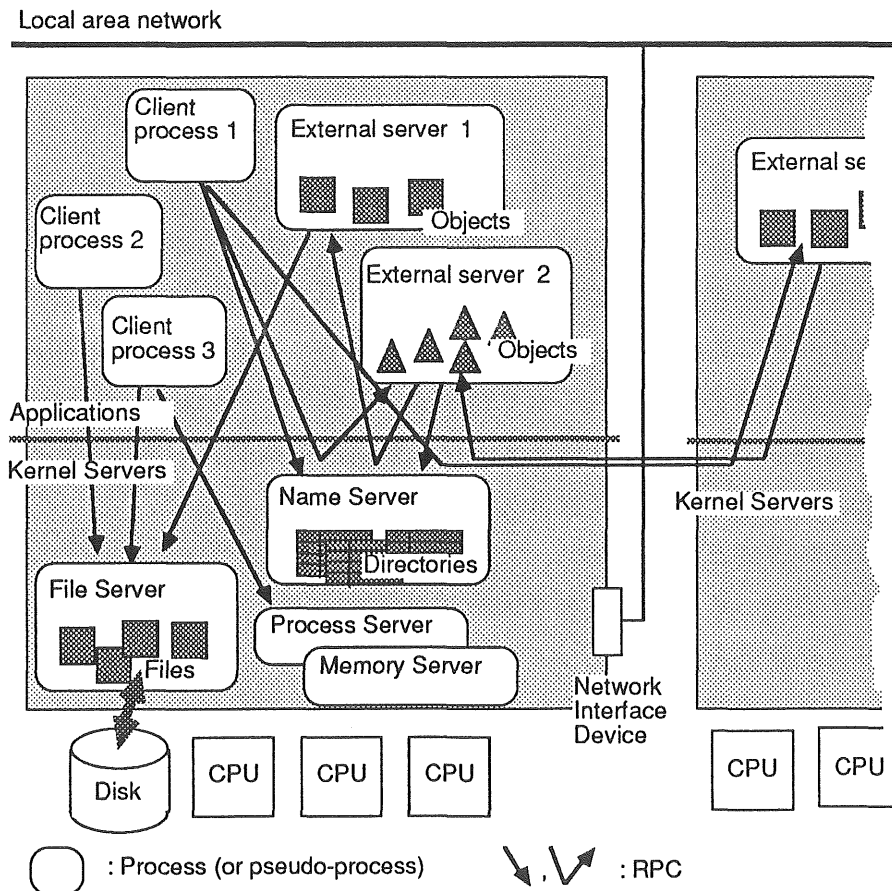


図2-2 カーネル・サーバと外部サーバ

Figure 2-1: Kernel servers and external servers.

## 2. 5 オブジェクトとサーバ

ReSCシステムは、object-basedシステムである。たとえば、ファイル、ディレクトリ、プロセスのようなソフトウェア的な資源は、オブジェクト (object) として統一的に扱われる。システムの機能は、そのようなオブジェクトを通じて提供される。

オブジェクトは、それぞれのサーバにより管理される。サーバとは、クライアント・プロセスに対してサービスを提供するプロセス、プロセスの集合、あるいは、カーネル中の疑似プロセスである。疑似プロセス (pseudo-process) とは、カーネル中に在り、プロセスではないが、クライアント・プロセスから普通のプロセスと同じように見えるものである (図2-2)。

サーバは、オブジェクトを管理するオブジェクト、すなわち、オブジェクトのメタ・オブジェクト (meta-object) に相当する。オブジェクトの操作は、そのオブジェクトのサーバに対してメッセージを送ることにより行われる。サーバ自身は、プロセス・オブジェクトであり、カーネル内のプロセス・サーバにより管理される。またカーネルは、プロセス間通信機能、すなわち、オブジェクトに対してメッセージを送る機能を提供する。プロセス間通信については、2. 7節において述べる。

カーネルは、サーバのサーバ、すなわち、メタ・オブジェクトのメタ・オブジェクトになっている。そしてカーネル自身は、ハードウェアにのみ依存し、ソフトウェア的には自立 (self-support)、および、自律 (autonomic) している。カーネルは、カーネル外のオブジェクトやサーバに依存していない。この性質により、カーネルの動作が安定し、カーネルの開発が容易になる。

## 2. 6 オブジェクト識別子

本システムでは、システム・レベルのオブジェクトの統一的な扱いを可能にするために、一様なオブジェクト識別子を導入した。このオブジェクト識別子の特徴は、位置依存のオブジェクトと独立独立のオブジェクトの両方を統一的に識別可能である点にある。

### 2. 6. 1 オブジェクト識別子の構造

本システムのオブジェクト識別子は、次の3つ組から構成される。

<サイト識別子、サーバ識別子、オブジェクト番号>

サイト識別子 (site-id) は、そのオブジェクトを管理するサーバが存在するサイトを示す数である。(サイト識別子については、2. 6. 2項において詳しく述べる。)

サーバ識別子 (server-id) は、そのオブジェクトを管理するサーバに割り当てられ

た、遠隔手続き呼出しのポート番号である。ネットワーク上に分散したプロセスの集合が、1つのサービスを提供する場合、それらのプロセスは、1つのサーバ識別子を共有する。オブジェクト番号 (object-number) は、各サーバにより管理される番号である。

## 2. 6. 2 サイト識別子

ReSCでは、位置依存のオブジェクトと位置独立のオブジェクトの両方を支援する。これを実現するために、次の3種類のサイト識別子を導入した。

### (1) サイト一意型サイト識別子 (a site-uniq site-id)

これは、各サイトに一意になるように割り当てた数である。インターネット [101] におけるホスト・アドレスに対応する。

### (2) 局所サイト型サイト識別子 (the local site site-id)

これは、局所サイトを示す数である。インターネットにおける、ソフトウェア・ループバック・ネットワーク・インタフェース (software loopback network interface) のアドレスに相当する。多くのインターネットのホストにおいて、このアドレスは、ホスト名 "localhost" により参照されるアドレス、すなわち、127.0.0.1 になっている。

### (3) 放送型、あるいは、マルチキャスト型サイト識別子

(a broadcast or multicast site-id)

これは、放送 (broadcast)、あるいは、マルチキャスト (multicast) を示す数である。インターネットにおける放送、または、マルチキャストを表すアドレスに相当する。

ここで、(1) を含むオブジェクト識別子をサイト一意型オブジェクト識別子 (a site-uniq object-id) とよぶことにする。同様に、(2) 含むオブジェクト識別子を局所サイト型オブジェクト識別子 (a local-site object-id)、(3) を含むオブジェクト識別子を、放送型オブジェクト識別子 (a broadcast object-id)、あるいは、マルチキャスト型オブジェクト識別子 (a multicast object-id) とよぶことにする。

サイト一意型オブジェクト識別子を利用することにより、位置依存のオブジェクトを識別することができる。したがって、サイト一意型オブジェクト識別子を、位置依存型オブジェクト識別子 (a location-dependent object id) とよぶ。以下で述べるように、局所サイト型、放送型、あるいは、マルチキャスト型のオブジェクト識別子を用いることで、位置独立なオブジェクトを識別することができる。このようなオブジェクト識別子を、位置独立型オブジェクト識別子 (a location-independent



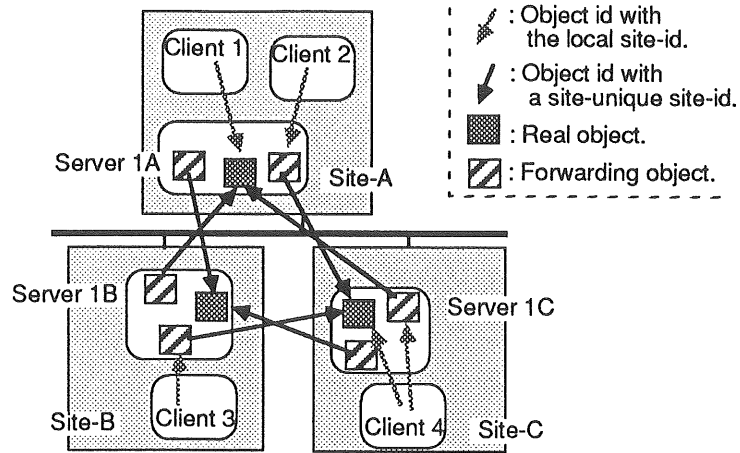


図 2 - 3 局所サイト型オブジェクト識別子による  
位置独立のオブジェクトの識別

Figure 2-3: Identifying location independent objects  
by the object-id with the local site-id.

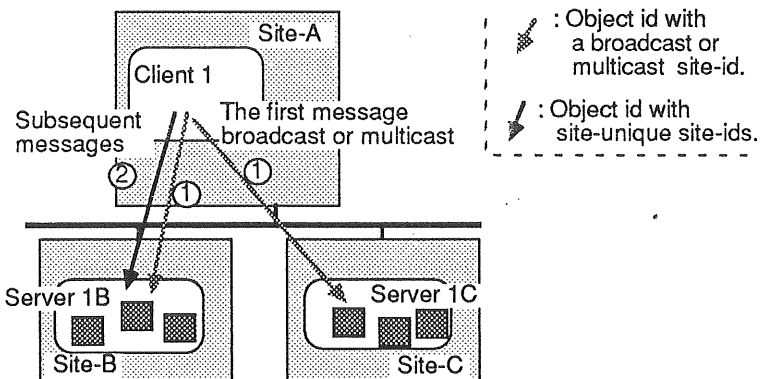


図 2 - 4 放送型、マルチキャスト型オブジェクト識別子による  
位置独立のオブジェクトの識別

Figure 2-4: Identifying a location independent object  
by the object id with a broadcast or multicast site-id

object-id) とよぶことにする。

局所サイト型オブジェクト識別子により指し示されたオブジェクトを利用する場合、クライアントは、自分自身と同一のサイト上のサーバに要求を送る（図 2 - 3）。そのサーバは、操作対象となるオブジェクトの位置を知っている。自分自身がそのオブ

ジェクトを管理している場合、サーバは、その要求に応答する。そうでない場合、別のサイト上に在り、実際にそのオブジェクトを管理しているサーバに、その要求メッセージを転送 (forward) する。

放送型、あるいは、マルチキャスト型オブジェクト識別子により指し示されたオブジェクトを利用する場合、クライアントは、放送、または、マルチキャストを使って要求を送る (図 2-4)。以後、応答を返したサーバと通信する。これは、Vシステムや Amoeba システムなど他の分散型オペレーティング・システムで用いられている方法とよく似ている。相違点は、本システムの場合、オブジェクト識別子中のサーバ識別子を見ただけでは、サーバの位置を特定する事ができないこと、および、複数のサーバが応答する可能性があることである。

### 2. 6. 3 位置情報を含むオブジェクト識別子に関する議論

オブジェクト識別子に位置情報を含めることについては、議論がある。多くの分散型オペレーティング・システムでは、位置情報を含まないオブジェクト識別子が好んで使われている。その理由は、システムのレベルでオブジェクトの複製やオブジェクトの移動を実現することが可能になるからである。

ReSCでは、以下の理由によりオブジェクト識別子に位置情報を含めることにした。

- (1) 分散型システムには、位置独立のオブジェクトに加えて、位置依存のオブジェクトが存在する。位置依存のオブジェクトの例としては、物理ディスク装置、ディスプレイ、キーボード、フロッピ・ディスク装置があげられる。
- (2) 並列/分散応用プログラムは、位置依存のオブジェクトを利用する。たとえば、分散応用プログラムの1つとして、複製を行うファイル・サーバのプログラムを書くことを考える。この時、複製ファイル・サーバは、オブジェクトの位置に注意を払う必要がある。もし1つのファイルについて、全てのコピーが1つのサイトに置かれたならば、複製は働かない。(オブジェクトの堆積による複製ファイル・サーバについては、5. 6. 1項で述べる。)
- (3) 位置独立のオブジェクトは、2. 6. 2項で述べた位置独立型オブジェクト識別子(局所サイト、放送、あるいは、マルチキャストを意味するサイト識別子を含むオブジェクト識別子)を使って識別することができる。

文献 [47] では、オブジェクト (文献 [47] では、ファイルのみ) の名前付け (名前には、識別子を含む) とオブジェクトの位置に関して、次のような概念を規定している。

- ・位置透明性 (location transparency)

オブジェクトの名前には、それが物理的にどこに置かれているかを示す一切の示唆が含まれない。

- ・位置独立性 (location independence)

オブジェクトの物理的な位置が変わっても、オブジェクトの名前を変更する必要がない。

ここで、位置独立性の方が強い概念である。すなわち、位置独立性を実現するためには、位置透明性に加えて、オブジェクトの位置が変化した前後 (オブジェクトの移動の前後) で、アクセスの透明性を保つために、動的にオブジェクトの位置を特定する機能が必要になるからである。

R e S C のオブジェクト識別子に対しては、この分類方法は、無効である。すなわち、2. 6. 2 項で述べた位置依存型オブジェクト識別子 (サイト一意型オブジェクト識別子) は、位置透明性も位置独立性も満たしていない。しかしながら、位置独立型オブジェクト識別子 (局所サイト型、放送型、マルチキャスト型) は、位置独立性を満たしている。このように、R e S C のオブジェクト識別子は、性質が異なるオブジェクト識別子が同一のドメインに含まれている点に特徴がある。

## 2. 7 遠隔手続き呼出し

カーネルは、プロセス間通信プリミティブとして遠隔手続き呼出し (RPC, Remote Procedure Call) を提供する。遠隔手続き呼出しは、同期式の通信を実現する。したがって、カーネルにおいてパラメタのバッファリングを行う必要がなく、データのコピーのオーバーヘッドを軽減し、効率を改善することができる。また、遠隔手続き呼出しは、通常の手続き呼出しと同じ様に扱うことができるので、使いやすい [98]。また、非同期式の通信が必要となる局面では、第3章で述べる軽量プロセスを用いることにより、それを実現することも可能である。さらに、エラーを考えると、ほとんど全てのサーバの呼出しは、必ず結果を返すので、単方向の通信プリミティブと比較して双方向の通信プリミティブである遠隔手続き呼出しが適している。

遠隔手続き呼出しは、A m o e b a、V システムをはじめとして、他の分散型オペレーティング・システムにおいて広く採用されている [98] [12] [54]。

### 2. 7. 1 位置一様性とサーバ型一様性

R e S C の遠隔手続き呼出しは、次の2つの一様性を備えている。

(1) 位置一様性：

クライアントは、呼び出すサーバが局所サイト上にあるか、遠隔サイト上にあるか、サーバの位置によらず、同一の方法で呼び出すことができる。このことは、ある種の分散透明性を実現していることになる。ゆえに、

ReSCは、ネットワーク・オペレーティング・システムではなく、分散型オペレーティング・システムの範疇に分類される。

(2) サーバ型一様性：

クライアントは、呼び出すサーバが外部サーバであるか、カーネル・サーバであるか、サーバの型によらず、同一の方法で呼び出すことができる。本システムでは、カーネルの内部にプロセス・サーバやファイル・サーバを含んでいる。そして、カーネルに対するサービスの要求も、他のカーネル外のサーバの呼出しと同一のインターフェースにより行うことが可能になっている。この性質を実現することにより、システムの拡張性が高められる。その理由は、カーネル・サーバもカーネル外のサーバもクライアントから区別ができないため、高い機能を持つサーバをカーネル外に構築することが容易になるからである。

## 2. 7. 2 分散システムの多重視点

位置一様な遠隔手続き呼出しと位置が見えるオブジェクト識別子により、分散システムの多重視点 (multiple view) が実現される。並列／分散応用プログラムは、システムを分散システムとして見て、遠隔サイト上のサーバに遠隔手続き呼出しを行い、遠隔サイトにファイルやプロセスを生成することができる。同時に、逐次応用プログラムは、システムを集中型システムとして見て、遠隔サイト上のファイルを局所サイト上のファイルと同じ方法でアクセスすることができる。

図2-5(a)に、システムを分散型システムとして見る例を示す。この図において、Process 1 が Site-B 上のプロセス・サーバを遠隔手続き呼出しにより呼び出し、Site-B 上に Process 2 を生成している。Process 2 は、局所サイト上のファイル・サーバを呼び出し、ファイルを読み込んでいる。そして、フィルタ操作を行い、ヒットしたデータのみを Site-A 上の Process 1 に送っている。この結果、ネットワーク上を転送されるデータ量が削減され、Process 1 と Process 2 の間で並列処理が実現されている。

図2-5(b)に、システムを集中型システムとして見る例を示す。この図において、Process 1 は、局所サイト Site-A 上のファイルと全く同じ方法で遠隔サイト Site-B 上のファイルをアクセスしている。

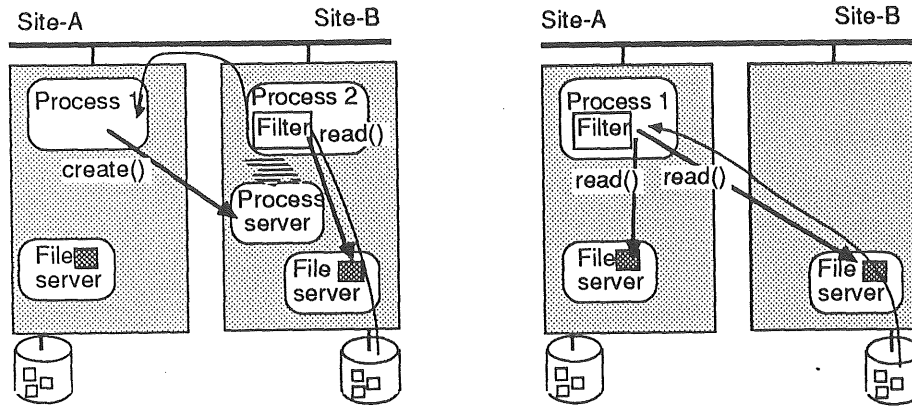


図 2-5(a) 分散型システムとして見る

図 2-5(b) 集中型システムとして見る

Figure 2-5(a): Viewing as a distributed system.

Figure 2-5(b): Viewing as a centralized system.

### 図 2-5 分散型システムの多重視点

Figure 2-5: Multiple views of the distributed system.

#### 2. 7. 3 遠隔手続き呼出しを行うカーネル・コール

遠隔手続き呼出しの機能は、次の5つのカーネル・コールにより提供される。(1)は、クライアントにより、残りは、サーバにより利用される。

(1) `rpc_send()`

これは、クライアントがサーバを呼び出すカーネル・コールである。クライアントは、サーバからの応答が到着するまでブロックされる。

(2) `rpc_receive()`

これは、サーバが、クライアントから要求を受け付けるカーネル・コールである。サーバは、要求が到着するまで、ブロックされる。

(3) `rpc_reply()`

これは、サーバが、クライアントへ応答を返すカーネル・コールである。

(4) `rpc_forward()`

これは、サーバが、他のサーバへ要求を転送するカーネル・コールである。

(5) `rpc_send_with_uid()`

これは、サーバが、別の利用者識別子を伴い他のサーバを呼び出すカーネル・コールである。

付録Aに、これらのカーネル・コールの詳細と、それらを利用したサーバとクライアントの例を示す。

サーバでは、普通、1つの遠隔手続き呼出しのセッションについて1つの軽量プロセスが生成される。その軽量プロセスは、クライアントからの要求を処理し、クライアントへ応答を行い、終了する。軽量プロセスを用いたサーバの構築については、5.9.1項において述べる。これに対して、軽量プロセスを生成することなく要求に対する処理を行い、応答を返す方法がある。この場合、Ada [68]におけるランデブ (rendezvous) と同じ動きになる。

## 2.8 認証

本システムは、マルチユーザのシステムである。したがって、一般の分散型オペレーティング・システムと同様に、利用者の認証 (authentication) を行う必要がある。本システムでは、LANという機密性が保たれた環境を利用して、集中型システムにおいて用いられている方法と同じように、暗号化の技術を用いない単純な方法により利用者の認証を行う。

### 2.8.1 利用者識別子

本システムでは、オブジェクトへのアクセス権の確認、セキュリティのチェックのために、統一的な利用者識別子を導入する。これは、次の組から構成される。

#### <利用者番号、マスク>

利用者番号は、利用者に対して一意になるように割り当てた整数である。マスクとは、対応する整数の有効範囲を現すものである。マスクのビットが1であることは、対応する整数のそのビットが有効であることを意味する。アクセス権の確認は、マスクのビットが1の部分について行われる。一般の利用者のマスクとしては、全ビット1 (32ビットの整数として、0xFFFFFFFF) の整数を用いる。サーバ・プロセスを起動したり、システム管理を行うための特権利用者 (UNIXにおけるスーパーユーザ) を示す利用者識別子としては、マスクが0である利用者識別子を用いる。これにより、特権利用者は、アクセス権のチェックをまったく受けることなく、全てのオブジェクトを操作することが可能になる。マスクを適当に設定することにより、インターネット [101] におけるサブネットと同様に、グループを設定することも可能である。付録A-3と付録A-4に、この利用者識別子の利用例を示す。

### 2.8.2 遠隔手続き呼出しとアクセス権の確認

カーネルは、個々の遠隔手続き呼出しにおいて、クライアント・プロセスに対応している利用者の識別子をサーバに渡す。サーバは、これを利用して、オブジェクトに対するアクセス権を確認する。

アクセス権を確認するために、プロセスの利用者の情報を付加することは、カーネルにおいて行う必要がある。たとえば、SunRPC [69] では、ライブラリにお

いて (UNIX の) 利用者識別子を遠隔手続き呼出しのヘッダに付加している。このため、容易に利用者識別子を偽ることが可能になっている。たとえば、次のような関数をリンクするだけで、利用者識別子が 0 の利用者 (UNIX では、スーパーユーザ) になりすますことができる。

```
getuid()
{
    return( 0 );
}
```

これは、SunRPC のライブラリの内部の、遠隔手続き呼出しのヘッダを作成する部分で、getuid() という関数を呼び出しているからである。getuid() は、本来、トラップ命令を含むシステム・コールである。この例では、利用者によりシステム・コールがオーバーライド (override) されている。その結果、アクセス権の確認が無効になっている。

最近の SunRPC では、特権利用者の確認を、遠隔手続き呼出しのヘッダ内の利用者識別子ではなく、インターネットの特権ポートを利用しているかどうかにより判定している。これにより、上記の問題は、解決されている。しかしながら、特権利用者以外の一般利用者については、問題が残されている。ResC システムでは、カーネルにおいて利用者の情報を付加するので、このような問題は生じない。

## 2. 9 名前サーバによる位置独立な名前付け

ResC では、オブジェクトの文字列の名前は、そのオブジェクトのサーバではなく、名前サーバ (name server) により管理される。名前サーバは、ディレクトリ (directory)、すなわち、オブジェクトの文字列の名前とオブジェクト識別子の組のリストを管理する。これにより、名前サーバは、ディレクトリ・サーバ (directory server) ともよばれる。名前サーバは、クライアントから検索 (lookup) 要求を受け付けると、渡された文字列が指定されたディレクトリに含まれているかどうかを検索する。その結果、含まれていた場合、対応するオブジェクト識別子をクライアントに返す。

2. 6 節において述べたように、オブジェクト識別子の中には、位置情報が含まれている。それにもかかわらず、文字列による名前付けでは、2. 6. 3 項で述べた位置透明性と位置独立性が実現されている。すなわち、オブジェクトの文字列の名前は、オブジェクトの位置が変化したとしても、変更する必要はない。このため、利用者は、局所サイト上のオブジェクトと同じ方法で遠隔サイト上のオブジェクトを扱うことができる。このような位置透明性のことを、文献 [60] では、名前透明性 (name tra

nsparency) とよんでいる。

ReSCの名前サーバ(ディレクトリ・サーバ)は、Amoebaのディレクトリ・サーバと類似している。Amoebaでは、ディレクトリ・サーバを利用して、書き換えがない(immutable)ファイルに対してアトミック・トランザクションを実現している[54]。ReSCのディレクトリ・サーバも、同じ目的に利用することが可能である。これに加えて、ReSCでは、書き換えがないオブジェクトに対して、オブジェクトの移動(object migration)を実現することができる。オブジェクトの移動は、その操作の実行後でオブジェクトの位置が変るようなアトミック・トランザクションにはほかならない。5.6.2項では、オブジェクトの堆積による書き換えがある(mutable)オブジェクトに対するオブジェクトの移動について述べる。

## 2.10 カーネル・サーバ群

この節では、カーネル・サーバの機能の概観を示す。カーネル・サーバには、プロセス・サーバ、名前サーバ、メモリ・サーバ、ファイル・サーバ等がある。カーネル・サーバは、下位層の並列/分散ハードウェアの機能を、過度に仮想化することなく、ほとんどそのままの形で提供する。そして、キャッシング、マッピングといった、並列/分散応用プログラムが好んで制御するような方針を、並列/分散応用プログラムに任せる。逐次応用プログラムに対する高度な機能は、外部サーバにおいて提供する。外部サーバについては、第5章で詳しく述べる。

カーネル・サーバは、プロセスではなく、カーネル内のモジュールである。それらのモジュールは、互に依存関係にある。しかしながら、カーネルは、カーネル外のモジュールに依存することはない。すなわち、カーネルは、自立して動作することが可能である。これにより、カーネルの動作が安定し、かつ、セキュリティを高めることが可能となる。

主なカーネル・サーバの機能を以下に示す。これらのカーネル・サーバのインタフェースを、付録Bに示す。

### 2.10.1 プロセス・サーバ

プロセス・サーバは、次の様な機能を提供する。

- (1) プロセスの生成・消滅
- (2) 仮想プロセッサの生成・消滅
- (3) プロセッサのスケジューリング
- (4) プロセッサの利用状況の提供

仮想プロセッサに関する機能の詳細は、第3章において述べる。プロセッサのスケジューリングは、通常の時分割システム(TSS, Time Sharing System)と同様な方針により行われる。(4)の機能は、第4章で述べるマッピング・コントローラにより



利用される。

#### 2. 10. 2 メモリ・サーバ

メモリ・サーバは、主記憶の管理を行う。主な機能を、以下に示す。

- (1) プロセスに対してメモリ・オブジェクトを割り当てる。
- (2) プロセスのメモリ・オブジェクトのマッピングを変更する。
- (3) カーネル内で用いるメモリを管理する。

#### 2. 10. 3 名前サーバ

名前サーバ（ネーム・サーバ、ディレクトリ・サーバ）は、2. 9 節で述べたように、与えられた文字列を対応するオブジェクト識別子に変換するサーバである。カーネル内の名前サーバは、基本的な機能を提供する。キャッシングや複製といった高度な機能を提供する名前サーバは、カーネルの外に外部サーバとして構築される。

#### 2. 10. 4 ファイル・サーバ

カーネル・ファイル・サーバは、標準ファイルに対するブロックごとの入出力機能を提供する。これは、NFSサーバ [59] の、ファイル・サーバとしての機能と同等の機能を提供する。（NFSサーバは、ディレクトリ・サーバとしての機能も持っている。）提供するファイルのモデルは、ランダム・アクセス・ファイルである。1 回の遠隔手続き呼出しごとに入出力を開始する点と長さを指定する。ファイルを開く操作は、存在せず、クライアントとの間に結合（connection）は作られない。ディスク・ブロックのキャッシングは、行われない。

カーネル・ファイル・サーバは、外部サーバにより提供可能な機能や、並列／分散応用プログラムから利用されない機能は、提供しない。たとえば、キャッシングやマッピングは、外部サーバにおいて提供することが可能である。さらに、並列／分散応用プログラムが、自分自身でキャッシングやマッピングを制御することで性能改善を目指す可能性がある。したがって、カーネル・ファイル・サーバは、キャッシングとマッピングの機能を提供しない。

#### 2. 10. 5 タイム・サーバ

タイム・サーバは、日付と時刻の情報を管理する。これは、応用プログラムが実行の統計を取る支援を行う。

### 2. 11 関連した研究

#### 2. 11. 1 カーネル化されたカーネル

オペレーティング・システムの構成法の 1 つに、カーネル化されたカーネル

(kernelized kernel) を使う方法がある [43]。この方法では、カーネルは、プロセス間通信と物理入出力だけを提供する。オペレーティング・システムの機能の大部分は、カーネル外のサーバ・プロセスにより提供される。これに対して、カーネル外のサーバを用いることなく、1つの巨大なカーネルを利用する方法がある。このようなシステムは、単層システム (monolithic system)、そのカーネルは、単層カーネル (monolithic kernel) とよばれる。

カーネル化されたカーネルは、普通、内部にデバイスを操作する部分だけを含んでいる。この部分は、カーネル外に実現するサーバによってのみ利用され、利用者から直接利用されることはない。利用者は、カーネル外のサーバを通じて間接的にカーネル内のデバイスを操作する部分を利用することになる。

カーネル化の目的は、オペレーティング・システムのできるだけ多くの部分をカーネル外のサーバとして利用者プログラムと同じように開発することを可能にすることである。その結果、機能の拡張と変更が容易になるという利点が得られる。また、オペレーティング・システム内部の並列性が、複数のサーバの並列実行という形で、容易に抽出される。さらに、単層カーネルを用いる方法と比較して、モジュール化が明確になるという利点もある。カーネル化の問題点としては、プロセス間の保護の壁を交差するオーバーヘッドが大きくなるため、効率が低下することである。

カーネル化されたカーネルは、しばしば分散カーネルとして実装される。分散カーネル (distributed kernel) とは、カーネルがネットワーク通信機能を提供し、かつ、ネットワーク上の全ての計算機においてカーネルのコピーが動作するようなものである。カーネル化された分散カーネルを利用して分散型オペレーティング・システムを構築している例としては、Amoebaシステム [54]、Vシステム [12]、Chorusシステム [1] があげられる。

本システムのカーネルの特徴は、並列／分散応用プログラムには単層カーネル、逐次応用プログラムには、カーネル化されたカーネルになっている点にある

(図1-3)。他のカーネルと異なり、本システムのカーネルは、内部に完全なサーバ (応用プログラムが直接呼び出すことができるサーバ) を含んでいる。これにより、並列／分散応用プログラムには、単層カーネルと同等の効率のよいサービスを提供することが可能となっている。

一方、逐次応用プログラムは、高度な機能を提供するカーネル外のサーバ (外部サーバ) を利用する。このことは、他のカーネル化されたカーネルを利用するシステムとほとんど同じである。5章では、オブジェクトの堆積というモデル化の手法を用いて、カーネル外の複数のサーバを統合して利用する方法について述べる。

## 2. 11. 2 マイクロカーネル

マイクロカーネル (micro-kernel) という言葉は、2. 11. 1項で述べたカーネル化されたカーネルを意味する言葉としても利用される。その他に、オペレーティング・システムのオペレーティング・システム、すなわち、メタ・オペレーティング・

システムを意味する言葉としても利用される [1] [2] [26]。メタ・オペレーティング・システムとは、仮想計算機 (VM、virtual machine) と同様に、オペレーティング・システムを実現するための環境を提供するものである。仮想計算機の場合、提供するモデルは、計算機のハードウェアである。そして、1台の物理的な計算機上に、複数の仮想計算機が動作し、それぞれの仮想計算機において独立したオペレーティング・システムが動作する。これに対して、マイクロカーネルの場合、提供するモデルが計算機のハードウェアではなく、ソフトウェア・アーキテクチャである。

メタ・オペレーティング・システムとしてのマイクロカーネルを利用した例としては、Machシステムがあげられる [2]。Machでは、マイクロカーネルにおいて、タスク、スレッド、ポート、メッセージ、メモリ・オブジェクトという5つの抽象を提供する。そして、マイクロカーネル上に、サーバとしてオペレーティング・システムが構築される。現在、UNIX、MS-DOS、MachintoshOS等がマイクロカーネル上のサーバとして動作している。

本システムのカーネルは、メタ・オペレーティング・システムとしてのマイクロカーネルではなく、それ自身として自律したオペレーティング・システムである。したがって、Machマイクロカーネル上で、本システムのカーネルを動作させることも考えられる。しかしながら、本システムのカーネルの機能は、Machマイクロカーネルの機能よりも低レベルの部分があるため、それを行うことは困難である。たとえば、3章で述べる仮想プロセッサの機能を、Machマイクロカーネルのスレッドやメモリ・オブジェクトの機能を利用して実現することは、困難である。

## 2. 12 まとめ

この章では、本ResCシステムのハードウェア構成、対象とする応用プログラム、全体構成、および、カーネルの機能について述べた。本システムの特徴は、カーネルの機能の設定にある。カーネルは、カーネル化されたカーネルの機能の他に、多重プログラミング・システムとしての資源割当ての方針を決定する部分や、並列/分散応用プログラムのための効率的なサービスを提供するサーバを含んでいる。また、システム・レベルのオブジェクト識別子は、位置依存のオブジェクトと位置独立のオブジェクトを同一のドメインで扱うことができるようになっている。位置依存のオブジェクト識別子と位置独立の遠隔手続き呼出しにより、システムは、分散型システムとも集中型システムとも見ることができるようになっている。

## 第3章 軽量プロセス

この章では、本 R e S C システムにおける軽量プロセスの機能と実現方式について述べる。R e S C では、マイクロプロセスと仮想プロセッサという概念を用いて軽量プロセスを実現する。ここでは、まず軽量プロセスの概念について整理を行い、従来の代表的な軽量プロセスの実現方式であるカーネル制御方式とコルーチン方式について説明する。次に、本マイクロプロセス／仮想プロセッサ方式について述べ、本方式の特徴である、応用固有の軽量プロセス間の同期・通信プリミティブとスケジューラの開発の支援について述べる。続いて、仮想プロセッサを提供するカーネルの構成方式、裸の共有メモリ型マルチプロセッサ上への実現、既存のシステム上でのカーネル機能のエミュレートについて述べる。最後に、実現したシステムの性能について述べ、関連した研究との比較を行う。なお、この章では、軽量プロセスを利用して並列応用プログラムを記述するプログラマに対して提供するシステムの機能について述べる。したがって、この章で利用者とは、並列応用プログラムのプログラマを、利用者レベルとは、カーネルではなく応用プログラムのレベルを意味するものとする。

### 3.1 軽量プロセスの概念

多重プログラミングのオペレーティング・システムにおいて、プロセスという言葉は、資源割当て (resource allocation)、および、保護 (protection) の単位を指すものとして用いられる。オペレーティング・システムは、プロセスごとに資源の管理表のエントリを設け、メモリ資源や CPU 資源の公平な分配を行う。さらに、各プロセスに別々のアドレス空間を割り当てることで、他のプロセスの誤動作、あるいは、悪意をもつ利用者による攻撃からプロセスを保護する。一方、並列処理の分野では、プロセスという言葉は、並列処理の単位を指すものとして用いられる。この場合、プロセスは、CPU 割当ての単位であり、CPU 以外の資源割当てや保護の単位ではない。

軽量プロセス (lightweight process) は、多重プログラミング・システムにおいて従来のプロセスとは異なる、1 応用プログラム内部の並列処理の単位としてのプロセスとして定義することができる。これに対して、従来のプロセスを重量プロセス (heavyweight process) とよぶ。たとえば、共有メモリ型マルチプロセッサにおいて並列処理を行う場合、軽量プロセスを並列処理の単位として用いることにより、CPU 処理と CPU 処理の重ね合せを、従来のプロセスを用いる方法よりも効率よく実現することが可能となる。また、軽量プロセスは、ウィンドウ・サーバやファイル・サーバなど、複数のクライアントを同時に扱う必要があるサーバを実現する際に有

効であることが知られている [ 9 7 ] 。この場合、クライアントごとに軽量プロセスを生成することにより、複数のサーバ・クライアント間の通信処理、入出力処理、および、CPU処理の重ね合せが実現される。

軽量プロセスは、その名が示すように、従来の資源割当てと保護の単位としてのプロセスと比較して、効率よく実現されることが求められる。具体的には、次のような操作の効率がよいことが求められる。

- ( 1 ) プロセスの生成、および、消滅
- ( 2 ) プロセス間の同期、および、通信
- ( 3 ) コンテキスト切替え

### 3. 2 軽量プロセス機能の実現目標

軽量プロセス機能の実現における目標は、次のようにまとめられる。

- ( 1 ) 実行時の効率がよい軽量プロセスを実現する

プロセスの生成・消滅、プロセス間の同期・通信、コンテキスト切替えのオーバーヘッドを小さくする。

- ( 2 ) 並列処理を実現する

CPU処理とCPU処理の重ね合せ、および、複数の通信処理、複数の入出力処理とCPU処理の重ね合せを実現する。

- ( 3 ) 応用固有の同期・通信プリミティブの実現を支援する。

並列応用プログラムでは、様々な種類の軽量プロセス間の同期・通信プリミティブが使われる。その同期・通信プリミティブとしては、セマフォやモニタ等、一般的なものが使われることもあるが、各応用プログラムに適した固有のものが使われることが非常に多い。このため、並列応用プログラムの開発においては、応用固有の同期・通信プリミティブを作ることが重要な作業となる。したがって、システムとしてこれをいかに支援するかが重要となる。

- ( 4 ) 応用固有のスケジューリングを実現する環境を提供する。

共有メモリ型マルチプロセッサにおいて並列処理を行う場合、応用固有のスケジューリングを行うことにより、著しい性能向上が可能な応用プログラムが存在する。たとえば、実行に先立ってスケジューリングを行うことにより、同期処理の回数を減らしたり、あるいは、並列性が高くなるように優先順位を決定することができる。このような応用プログラムに対しては、応用固有のスケジューリングを実現する環境を

提供することが重要となる。

(5) 軽量プロセス機能を実現するライブラリとカーネルの構成法を開発する。

オペレーティング・システムのカーネルは、カーネル上で動作する利用者プログラムと比較して、開発とデバッグが困難である。そのため、カーネルの構成法は、非常に重要である。2. 1 1. 1項で述べたように、オペレーティング・システムの機能の大部分をカーネルの外で動作するサーバ・プロセスにより提供することで、カーネルを小さくすることも試みられている。しかしながら、軽量プロセスを実現する上で、カーネルの外に移動することができない部分が存在する。したがって、そのようなカーネルを構成する方法が重要な研究課題となっている。

### 3. 3 代表的な軽量プロセス実現方式

軽量プロセスを実現する基本的な方法は、重量プロセスから保護の機能を省略することである。具体的には、アドレス空間を共有することである。こうすることにより、軽量プロセス間の通信を、軽量プロセスの共有変数を介して実現することが可能となり、効率が改善される。重量プロセスの場合、プロセス間通信によりプロセスが破壊されないようにオペレーティング・システムの保護機能が働く。

代表的な軽量プロセスの実現方式は、次の2種類に分類される。

#### (1) カーネル制御方式

この方式の基本的な考え方は、重量プロセスより保護の機能を省略することで、軽量プロセスを実現することである。軽量プロセス間の同期と通信は、同一のアドレス空間にあることを利用して、軽量プロセス間の共有変数を利用して実現される。軽量プロセスの生成・消滅、軽量プロセスのコンテキスト切替えは、カーネルにより制御される。この方式では、CPU処理の並列処理が可能である。この方式は、Machシステム[2][67]で用いられている。

#### (2) コルーチン方式

これは、軽量プロセスをコルーチンとして実現する方式である。この方式では、軽量プロセスの生成・消滅、軽量プロセス間の同期・通信など全ての操作は、利用者レベルにおいて実現される。この方式では、共有メモリ型マルチプロセッサにおいてCPU処理を並列に実行することができないが、CPU処理と入出力処理の重ね合わせについては、非同期入出力とソフトウェア割込みを用いて実現することは、可能である[27]。この方式は、SunOS[108]で用いられている。

以下では、マイクロプロセスと仮想プロセッサという概念を用いた軽量プロセスの

実現方式について述べる。本方式は、カーネル制御方式と比較して、効率的な軽量プロセスの実現が可能である。また、コルーチン方式と異なり、共有メモリ型マルチプロセッサにおいてCPU処理を並列に実行することができる。

### 3.4 マイクロプロセスと仮想プロセッサ

本システムでは、マイクロプロセスと仮想プロセッサという概念を用いて軽量プロセスを実現する。

マイクロプロセス (microprocess) は、並列処理の単位としてのプロセス、すなわち、軽量プロセスである。各マイクロプロセスは、重量プロセスの中にあり、独立のプログラム・カウンタとスタックを備えている。プロセスとマイクロプロセスの関係を図3-1に示す。マイクロプロセスの生成・消滅、マイクロプロセス間の同期・通信など、全ての制御は、利用者レベルにおいて行われる。オペレーティング・システムのカーネルは、一切介在しない。したがって、それらの制御が高速に実行される。

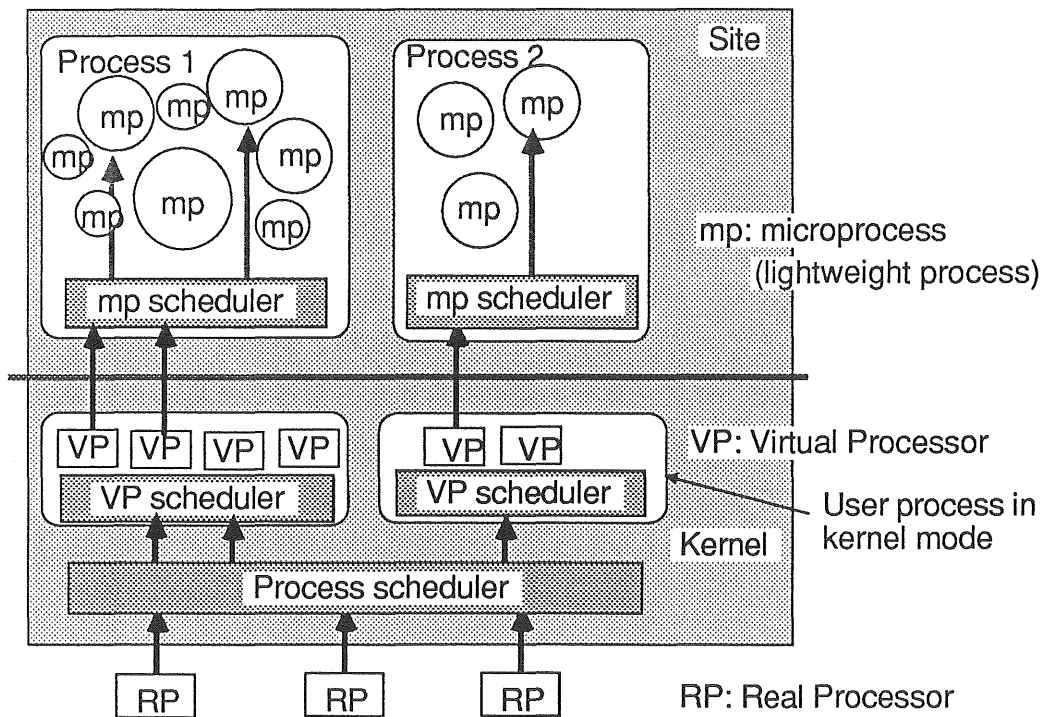


図3-1 プロセス、マイクロプロセス、仮想プロセッサの関係

Figure 3-1: Processes, microprocesses and virtual processors.

仮想プロセッサ (virtual processor) は、カーネルが利用者プロセスに対して実プロセッサ (real processor、ハードウェアのプロセッサ) を割り当てるためのエンタリである。図3-1に、プロセスと仮想プロセッサの関係を示す。各仮想プロセッサ

サは、1つの利用者プロセスと対応しており、その利用者プロセスに割り当てられた資源、利用者識別子、アクセス権、プロセスの優先順位などを共有している。

仮想プロセッサは、利用者レベルの物理的な並列処理を行うための仕組みである。共有メモリ型マルチプロセッサでは、カーネルは、1つのプロセスに複数の実プロセッサを割り当てる。各実プロセッサは、そのプロセスの仮想プロセッサの中から、実行可能なものを選択し、実行する。その結果、制御が利用者空間へ移動する。そして、その複数の仮想プロセッサにより複数のマイクロプロセスが並列に実行される。単一プロセッサでは、CPU処理を並列に実行することはできないが、複数の仮想プロセッサを生成することで、1つのCPU処理、複数の入出力処理、および、複数の通信処理を並列に実行することが可能となる。たとえば、利用者プロセスが、あるマイクロプロセスにおいて入出力処理を行うカーネル・コールを発行した場合、制御は、カーネル空間へ移動する。カーネルは、カーネル・コールを発行した仮想プロセッサを、入出力完了待ち状態に変え、同一プロセスの別の実行可能な仮想プロセッサを選択し、コンテキストを切り替える。この結果、制御は、再び利用者空間へ移動し、別の実行可能なマイクロプロセスが実行される。

プロセスが生成される時には、自動的に1個の仮想プロセッサが生成され、割り当てられる。これにより、逐次応用プログラムと並列応用プログラムの調和が図られる。逐次応用プログラムは、そのまま単一仮想プロセッサにより処理を進める。これにより、従来の逐次応用プログラムを仮想プロセッサに関して全く変更することなく共有メモリ型マルチプロセッサ上で利用可能となる。ゆえに、従来の単一プロセッサ・システムからの移行は、容易である。

一方、並列応用プログラムは、カーネル・コールを発行し、複数の仮想プロセッサを得る。並列応用プログラムからは、専用の共有メモリ型マルチプロセッサを使っているように見えるかもしれない。しかしながら、本システムは、多重プログラミング・システムであり、実プロセッサのスケジューリングは、カーネルにより行われる。すなわち、カーネル内のスケジューラが、各プロセスに割り当てる実プロセッサの数を調整する。このスケジューラをプロセス・スケジューラ (process scheduler)、あるいは、(1サイト内の) 大域スケジューラ (global scheduler) とよぶ

(図3-1)。そのサイトの負荷が低く、他に実プロセッサを必要としているプロセスが存在しない場合、全ての実プロセッサが1つのプロセスに割り当てられる。その結果、そのプロセスのプログラムは、資源の専有が可能な単一プログラミング・システムにおける実行と同等の速度で実行される。しかしながら、システムの負荷が高く、他に実プロセッサを要求しているプロセスが存在するならば、1つのプロセスに対して割り当てられる実プロセッサの数が減らされる。

プロセス、マイクロプロセス、仮想プロセッサについてまとめると次のようになる。

#### (1) プロセス

プロセスは、資源割当ての単位である。資源としては、メモリ、アドレス空



間、ファイルなどがある。実プロセッサ資源の割当ての単位もプロセスである。カーネル内のプロセス・スケジューラは、各プロセスに割り当てる実プロセッサの個数を決定する。

#### (2) マイクロプロセス

マイクロプロセスは、1 応用プログラム内の並列処理の単位である。マイクロプロセスは、独立したスタックとプログラム・カウンタを持つ。マイクロプロセスの生成・消去、コンテキスト切替え、マイクロプロセス間の通信など全ての制御は、応用プログラムのレベルで行われる。この結果、並列処理を行う際にオペレーティング・システムが介在する場面を極力減らすことが可能となり、効率が改善される。

#### (3) 仮想プロセッサ

仮想プロセッサは、カーネルがプロセスに実プロセッサを割り当てるためのエントリである。共有メモリ型マルチプロセッサでは、物理的並列処理単位となる。逐次プロセッサにおいても、複数の仮想プロセッサを生成することにより、1つのCPU処理、複数の入出力処理、および、複数のプロセス間通信処理を同時に行うことが可能となる。

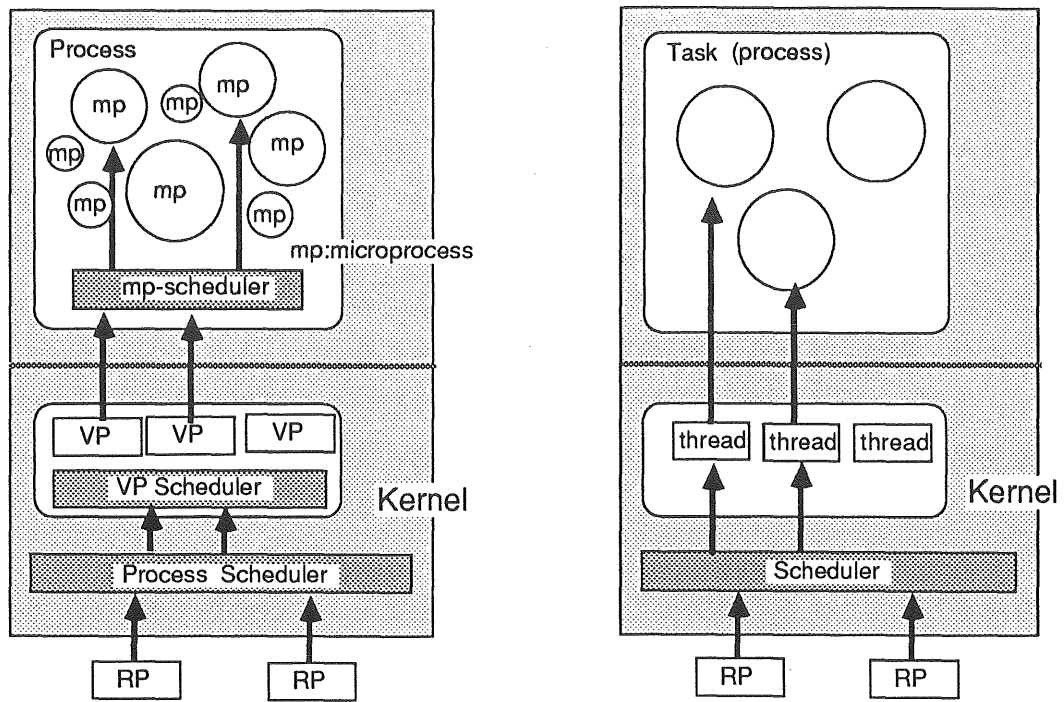
### 3. 4. 1 他の実現方式との比較

本方式、カーネル制御方式、および、コルーチン方式よる軽量プロセスの実現の様子を図3-2に示す。

#### ■カーネル制御方式との比較

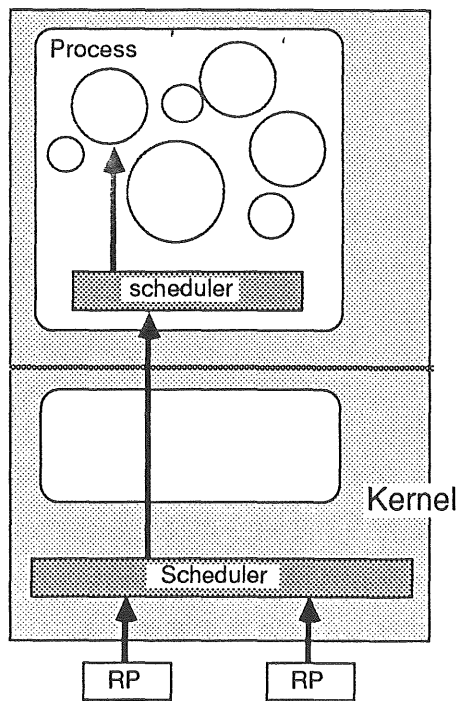
図3-2(a)に、本方式、図3-2(b)に、Machシステムにおけるカーネル制御方式による軽量プロセスを示す。Machでは、軽量プロセス(スレッド, thread)をカーネルが直接制御している。これに対して、本方式では、軽量プロセス(マイクロプロセス, mp)を応用プログラム・レベルのスケジューラが制御している。

本方式における仮想プロセッサは、Machシステムにおけるスレッドと対応することがわかる。しかしながら、本方式の仮想プロセッサの場合、利用者レベルで軽量プロセスが実現されていることを前提として機能が設計されている点が、大きく異なる。たとえば、本システムの仮想プロセッサは、1度生成されるとプロセス全体が終了するまで存在し続ける。また、本方式の仮想プロセッサには、仮想プロセッサごとの固有領域が存在し、利用者レベルの軽量プロセス(マイクロプロセス)を実現するライブラリにおいて利用される。仮想プロセッサの固有領域については、3.6.3項において詳しく述べる。



(a) マイクロプロセス/  
仮想プロセッサ方式

(b) カーネル制御方式 (Mach)



(c) コルーチン方式

VP:Virtual Processor  
 RP:Real Processor  
 ○ : Lightweight process  
 軽量プロセス

図 3 - 2 軽量プロセス実現方式の比較

Figure 3-3: Comparison of lightweight process implementation methods.

## ■ コルーチン方式との比較

図3-2(a)に示した本方式と図3-2(c)に示したコルーチン方式を比較すると、本方式におけるマイクロプロセスとコルーチンが非常に類似していることがわかる。マイクロプロセスとコルーチンの相違点は、コルーチンの場合、1度に1ルーチンしか実行されないが、マイクロプロセスの場合、複数の仮想プロセッサにより並列に実行されることである。そのため、共有メモリ型マルチプロセッサにおいて並列処理が可能となる。単一プロセッサにおいても、1つのCPU処理、複数の入出力処理、および、複数のプロセス間通信が並列に実行される。

マイクロプロセスの場合、軽量プロセス間の共有変数を参照する時にロックなどにより相互排除を行う必要がある。コルーチンの場合、その必要はない。また、コルーチンでは、あるルーチンでカーネル・コールを発行するとプロセス全体の実行が停止してしまう。マイクロプロセスでは、カーネル・コールを発行したマイクロプロセスとそれを実行していた仮想プロセッサは、一時停止するが、他の仮想プロセッサにより、プロセスの実行は継続される。

### 3.4.2 プロセッサのスケジューリング

本方式では、マイクロプロセス、プロセス、仮想プロセッサの3つのレベルにおいてプロセッサのスケジューリングを行う(図3-1)。それぞれ、次の様な機能をもつ。

#### (1) マイクロプロセス・スケジューラ

利用者プロセス中にあり、各応用プログラム固有の情報を利用して、応用固有のスケジューリングを実現する。コンテキスト切替えは、マイクロプロセスが同期プリミティブを発行した時に限り行われる。(横取りを行わない(non-preemptive)スケジューリングを行う。)

#### (2) プロセス・スケジューラ (process scheduler)

カーネル中にあり、多重プログラミング・システムにおける複数の利用者プロセス間の公平なプロセッサ資源の分配を行う。これを実現するために、量子時間(time quantum)ごとに、プロセッサの横取り(preemption)を行い、各利用者プロセスに対して割り当てる実プロセッサの数を調整する。

#### (3) 仮想プロセッサ・スケジューラ

カーネル中にあり、同一プロセス内の別の仮想プロセッサにコンテキストを切り替える。このコンテキスト切替えは、ある仮想プロセッサが入出力待ち、または、通信待ちのために実行を継続できなくなった時に限り行われる。(横取りを行わないスケジューリングを行う。)

このように機能を分割することにより、それぞれのスケジューラが単純化される。一方、スケジューラが多段になり、コンテキスト切替えのオーバーヘッドが大きくなることが予想される。しかしながら、スケジューラを実行する頻度がレベル間で異なり、圧倒的にマイクロプロセス・スケジューラの頻度が大きいので、仮想プロセッサ・スケジューラとプロセス・スケジューラは、性能に大きな影響を与えない。

本方式では、仮想プロセッサ・スケジューラにおいてプロセッサの横取りを行わない。この理由は、次の2点にある。

(1) マイクロプロセス・スケジューラを活用する。

仮想プロセッサ・スケジューラにおいて、プロセッサの横取りを行うならば、1つの並列応用プログラムの実行において、2つのスケジューラが同時に動作することになる。その場合には、利用者空間内のマイクロプロセス・スケジューラが機能しなくなる。

(2) 無駄なコンテキスト切替えを避ける。

本来、プロセッサの横取りは、複数のプロセス間において、プロセッサ資源の公平な利用を実現するための機能である。本方式では、資源割当ての単位は、プロセスであるので、同一プロセス内において資源の公平な利用を実現する必要がない。

### 3. 5 マイクロプロセスを実現するライブラリの機能と構造

本方式では、マイクロプロセスは、利用者空間内に実現されるので、カーネル・コールではなくライブラリにより支援される。このライブラリは、図3-3に示すように、層構造になっている。このため、各並列応用プログラムの要求に応じた部分的な利用が容易になっている。

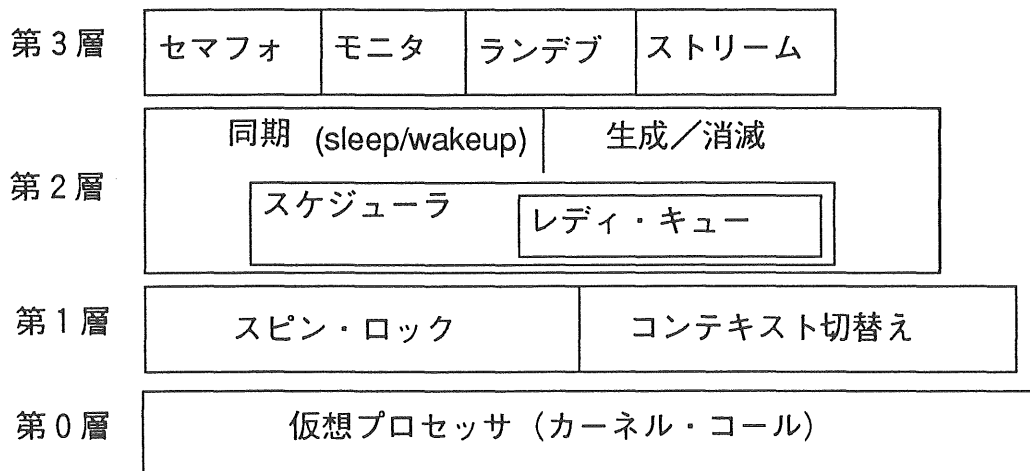


図3-3 マイクロプロセス・ライブラリの構造

Figure 3-3: The structure of the microprocess library.

第0層では、カーネルによって実現される仮想プロセッサに関連したカーネル・コールが提供される。この層を利用する並列応用プログラムは、仮想プロセッサだけを利用して、完全に独自の方式により軽量プロセスを実現することになる。仮想プロセッサに関連したカーネル・コールについては、3.6.1項で述べる。

第1層では、各プロセッサ・アーキテクチャやバス・アーキテクチャに独立なコンテキスト切替えとスピン・ロックを実現する機能が提供される。この層では、それらのアーキテクチャに独立したプログラムを記述することが可能となる。スピンロックは、test-and-set 命令や、swap命令のような、原始的な (atomic) 命令を用いて実現される。コンテキスト切替えは、主にレジスタの退避と回復により実現される。

第2層では、基本的なマイクロプロセスの生成・消滅機能と、sleep/wakeup モデルに基づく簡単な同期機能が提供される。このモデルの詳細については、3.5.1項において述べる。また、第2層は、スケジューラとレディ・キューを含んおり、これらは、応用固有のスケジューリングを行うために利用される。応用固有のスケジューリングについては、3.9節において述べる。

第3層では、セマフォ、モニタ、ランデブ、バイト・ストリーム (UNIXのパイプと同等の機能) 等の同期・通信プリミティブが提供される。これらのプリミティブのソース・コードは、各並列応用プログラム固有の同期・通信プリミティブを構築するプログラムにプロトタイプを与える。

並列応用プログラムは、通常、第2層の機能を用いて構築される。これは、セマフォやモニタを利用するよりも、低いレベルの機能を利用する方が各並列応用プログラム固有の同期・通信プリミティブを構築しやすいからである。

マイクロプロセス・ライブラリの主要な手続きの仕様を、付録Cに示す。

### 3.5.1 sleep/wakeup モデル

sleep/wakeup モデルでは、主に次の2つの手続きを用いて同期・通信プリミティブを構築する。

(1) mp\_sleep():

実行中状態 (running state) の自分自身 (マイクロプロセス) を、待ち状態 (blocked state) に移す。

(2) mp\_wakeup():

待ち状態にある (他の) マイクロプロセスを、実行可能状態 (runnable state、ready state) に移す。既に実行可能、あるいは、実行中の場合には、何もしない。

この他に、自分自身を実行可能状態にしたまま、他の実行可能なマイクロプロセスに制御を移すための手続き、mp\_yield()がある。各プロセッサは、実行可能状態にあるマイクロプロセスを選び、それを実行中状態に変え、そのマイクロプロセスへ制御を

移す。このモデルの特徴は、他の実行中のマイクロプロセスを待ち状態、あるいは、実行可能状態に移すプリミティブが存在しない点にある。

ここでは、典型的な同期プリミティブとして、2進セマフォを取り上げ、sleep/wakeup モデルについて述べる。図3-4(a)に、このモデルに基づく2進セマフォの実現例を示す。この図において、lock()、unlock() は、第1層の手続きであり、それぞれ、スピンロックによるロック／アンロックを行うものである。このプログラムでは、これらの手続きを用いて、際どい部分(critical section)の相互排除を行っている。

P命令の実現を図3-4(a)第11行に示す。第16行において、他のマイクロプロセスがそのセマフォを保持している(sem->inuse が TRUE)の時には、自分自身(current)を、待ち行列に保存し、セマフォの資源を開放する。そして、第19行において、mp\_sleep() を呼び出し、現在実行中のマイクロプロセスを待ち状態にする。セマフォが利用されていない時(sem->inuse が FALSE)には、自分自身が利用していることをマーク(第22行)している。

V命令の実現を、図3-4(a)第28行に示す。P命令の中で mp\_sleep() により実行を中断したマイクロプロセスは、第37行の mp\_wakeup() により、実行可能状態に戻される。

第15行において、現在実行中のマイクロプロセス(current)に掛けられたロックは、mp\_sleep() の内部で解除される。この解除を、この部分で先に行うことはできない。それは、sem->lock が開放された瞬間に、第36において、mp\_wakeup() が呼び出される可能性があるからである。実行可能中のマイクロプロセスに対して mp\_wakeup() が適用された場合、何も行われないため、そのシグナルが落ちる。

### 3.5.2 他の実現方式と比較

この項では、2進セマフォ同期プリミティブの実現を通じて、カーネル制御方式、および、コルーチン方式との比較を行う。

#### ■カーネル制御方式との比較

図3-4(b)に、カーネル制御方式の軽量プロセスにおけるセマフォの実現を示す。カーネル制御方式の例として、Machシステムのスレッドを取り上げた。ここで、3.5.1項と同様に、スピン・ロックによりロック機構を実現するライブラリ関数 lock()、unlock()、および、FIFOのキューを実現するライブラリ関数 enqueue()、dequeue() を利用する。

第28行、および、第47行に現れる mach\_msg() がMachシステムのカーネル・コールである。この実現では、メッセージ・パッシングを用いてスレッドの一時停止／実行再開を実現している。すなわち、スレッドの一時停止を行う時には、メッセージを受け取るカーネル・コール(mach\_msg() の第2引数が MACH\_RCV\_MSG (Mach receive message)) を発行している。逆に、一時停止したスレッドの実行再開を要求

```

1: extern private struct mpcb *current ;
2:
3: struct sem
4: {
5:     int          inuse ;
6:     lock_t       lock ;
7:     struct mpcb  *queue ;
8: };
9: typedef struct sem sem_t ;
10:
11: P( sem )
12:     sem_t *sem ;
13: {
14: loop:  lock( &sem->lock );
15:         lock( &current->lock );
16:         if( sem->inuse )
17:         {
18:             enqueue( current, &sem->queue );
19:             unlock( &sem->lock );
20:             mp_sleep( 0, 0 );
21:             goto loop;
22:         }
23:         sem->inuse = TRUE ;
24:         unlock( &sem->lock );
25:         unlock( &current->lock );
26: }
27:
28: V( sem )
29:     sem_t *sem ;
30: {
31:     struct mpcb *next ;
32:     struct mpcb *dequeue();
33:
34:     lock( &sem->lock );
35:     sem->inuse = FALSE ;
36:     next = dequeue( &sem->queue );
37:     if( next != MP_NULL )
38:         mp_wakeup( next );
39:     unlock( &sem->lock );
40: }

```

図 3 - 4 ( a ) マイクロプロセッサ・ライブラリによる 2 進セマフォの実現

図 3 - 4 2 進セマフォの実現

```

1: struct proc
2: {
3:     struct proc    *link ;
4:     mach_port_t    port ;
5:     mach_msg_header_t msg ;
6: };
7: typedef struct proc *proc_t ;
8: extern proc_t proc_self();
9:
10: struct sem
11: {
12:     int            inuse ;
13:     lock_t         lock ;
14:     proc_t         *queue ;
15: };
16: typedef struct sem sem_t ;
17:
18: P( sem )
19:     sem_t *sem ;
20: {
21: loop:  lock( &sem->lock );
22:         if( sem->inuse )
23:             {
24:                 mach_msg_header_t msg;
25:                 proc_t current = proc_self();
26:                 enqueue( current, &sem->queue );
27:                 unlock( &sem->lock );
28:                 mach_msg( &msg, MACH_RCV_MSG, 0, sizeof(msg),
29:                         current->port, MACH_MSG_TIMEOUT_NONE, MACH_PORT_NULL );
30:                 goto loop;
31:             }
32:         sem->inuse = TRUE ;
33:         unlock( &sem->lock );
34: }
35:
36: V( sem )
37:     sem_t *sem ;
38: {
39:     proc_t next ;
40:     extern proc_t dequeue();
41:     lock( &sem->lock );
42:     sem->inuse = FALSE ;
43:     next = dequeue( &sem->queue );
44:     unlock( &sem->lock );
45:     if( next != MACH_PORT_NULL )
46:         {
47:             mach_msg( &next->msg, MACH_SEND_MSG, sizeof(next->msg), 0,
48:                     MACH_PORT_NULL, MACH_MSG_TIMEOUT_NONE, MACH_PORT_NULL);
49:         }
50: }

```

図 3 - 4 ( b ) カーネル制御方式における 2 進セマフォの実現

図 3 - 4 2 進セマフォの実現



```

1: extern struct coroutine *current ;
2:
3: struct sem
4: {
5:     int          inuse ;
6:     struct coroutine *queue ;
7: };
8: typedef struct sem sem_t ;
9:
10:
11: P( sem )
12:     sem_t *sem ;
13: {
14: loop:   if( sem->inuse )
15:         {
16:             enqueue( current, &sem->queue );
17:             c_sleep();
18:             goto loop;
19:         }
20:     sem->inuse = FALSE ;
21: }
22:
23: V( sem )
24:     sem_t *sem ;
25: {
26:     struct coroutine *next ;
27:     struct coroutine *dequeue();
28:     sem->inuse = FALSE ;
29:     next = dequeue( &sem->queue );
30:     if( next != C_NULL )
31:         c_wakeup( next );
32: }

```

図3-4 (c) コルーチン方式による2進セマフォの実現

図3-4 2進セマフォの実現

する場合には、メッセージを送信するカーネル・コール (mach\_msg()) の第2引数が MACH\_SEND\_MSG (Mach send message) を発行している。この方法は、Mach 3.0 の C-Threads ライブラリにおいて用いられている。

Mach のスレッドを用いた実現 (図3-4 (b)) と、本システムにおける実現 (図3-4 (a)) を比較すると、Mach の場合に、カーネル・コールになっている部分が、本システムでは、(カーネル・コールを利用しない) ライブラリにより実現されている。また、コンテキスト切り替えの必要がない場合においては、カーネル制御方式の場合においても、カーネル・コールを発行しないことがわかる。

#### ■ コルーチン方式との比較

図3-4 (c) に、コルーチン方式の軽量プロセスにおけるバイナリ・セマフォの実現を示す。c\_sleep()、c\_wakeup() は、3.5.1項において述べた mp\_sleep(),

mp\_wakeup() と同等なライブラリ関数である。

コルーチン方式におけるセマフォの実現は、本方式における実現と非常に類似している。これは、本方式における実現からロックの制御を削除したものに相当する。

### 3. 5. 3 実行の軌跡の記録と解析

本マイクロプロセス・ライブラリには、マイクロプロセスの生成・消滅や、コンテキスト切替えが起きたことをイベントとしてその時刻を記録する機能がある。この機能は、並列応用プログラムを開発する上で、並列性の抽出の様子を確認するために利用される。

記録されたイベントを用いて、各マイクロプロセス、および、仮想プロセッサごとのCPU利用時間、入出力時間、コンテキスト切り替えの回数といったの統計情報を得ることができる。その例を、図3-5に示す。この例では、マイクロプロセスが9個(0番~8番)、仮想プロセッサが3個生成されている。マイクロプロセス関連の情報において、マイクロプロセス0番は、最初のマイクロプロセスを生成するスケジューラである。1番から8番は、普通のマイクロプロセスである。この部分の列は、それぞれ、次のような意味を持つ。

- (1) マイクロプロセスの番号
- (2) そのマイクロプロセスが利用したCPU時間の合計とその割合
- (3) CPUを利用した回数(CPUバーストの回数)
- (4) そのマイクロプロセスが行った入出力の時間の合計とその割合
- (5) 入出力回数
- (6) CPU時間と入出力時間の合計とその割合
- (7) コンテキスト切替えの回数

この例では、マイクロプロセス1番が全体の19%の実行時間を占めていること、特に入出力が重たくなっていること等が読み取れる。

仮想プロセッサ関連の部分では、次のような情報が得られる。

- (1) 仮想プロセッサ番号
- (2) CPU稼働時間(この仮想プロセッサがマイクロプロセスを実行していた時間)とその割合。
- (3) 入出力時間(この仮想プロセッサにより入出力行っていた時間)とその割合
- (4) その仮想プロセッサがアイドル状態にあった時間とその割合
- (5) この軌跡をファイルに保存するために要した時間と割合
- (6) その他の時間(時刻の入手、コンテキスト切替えなど)と割合
- (7) 合計とその割合

最終行は、全仮想プロセッサの平均である。図3-5の例では、プロセッサのアイドル時間が15%であったことがわかる。

```

## microprocess information
No  cputime  % / #ex  iotime  % / #io  total  % / #++
0   0.000  0 / 1   0.000  0 / 0   0.000  0 / 1
1   0.140  8 / 67  0.370  36 / 62  0.510  19 / 129
2   0.370  22 / 20  0.000  0 / 0   0.370  14 / 20
3   0.100  6 / 66  0.180  18 / 65  0.280  10 / 131
4   0.360  21 / 18  0.000  0 / 0   0.360  13 / 18
5   0.110  6 / 66  0.180  18 / 65  0.290  11 / 131
6   0.350  21 / 13  0.000  0 / 0   0.350  13 / 13
7   0.130  8 / 66  0.160  16 / 65  0.290  11 / 131
8   0.140  8 / 72  0.130  13 / 65  0.270  10 / 137
##   1.700 100 / 389  1.020 100 / 322  2.720 100 / 711

```

```

## virtual processor information
No   cpu  %|   io  %|  idle  %| trace  %| kernel  %| total  %|
0   0.600 54| 0.240 21| 0.230 21| 0.000 0| 0.050 4| 1.120 100|
1   0.560 50| 0.370 33| 0.160 14| 0.000 0| 0.030 3| 1.120 100|
2   0.540 48| 0.410 37| 0.130 12| 0.000 0| 0.040 4| 1.120 100|
##   1.700 51| 1.020 30| 0.520 15| 0.000 0| 0.120 4| 3.360 100|

```

図 3 - 5 実行の軌跡の記録より得られる統計情報

Figure 3-5: Execution statistics from microprocess traces.

記録されたイベントを、実行の軌跡として図 3 - 6 に示すように視覚化することも可能である。1 番から 8 番において、黒い部分がマイクロプロセスの実行を表す。i0, i1, i2 は、それぞれ仮想プロセッサの 0 番、1 番、2 番が、他に実行するマイクロプロセスが存在しないため、アイドル状態になったことを表す。並列応用プログラムの開発者は、このような情報を利用して、並列性の抽出の様子を観察することが可能である。

マイクロプロセス・ライブラリでは、具体的には、次のようなイベントを記録している。

- (1) 生成 (created)
- (2) 終了 (exit)
- (3) 中断 (sleep)
- (4) 実行可能化 (waked)
- (5) プロセッサの譲渡 (yield)
- (6) 実行開始 (scheduled)
- (7) 入出力開始 (io start)
- (8) 入出力終了 (io stop)
- (9) プロセス全体の終了 (process exit)
- (10) 仮想プロセッサの一時停止 (vp suspend)
- (11) 仮想プロセッサの実行再開 (vp resume)
- (12) 軌跡の記録のための入出力開始 (trace file io start)
- (13) 軌跡の記録のための入出力終了 (trace file io stop)

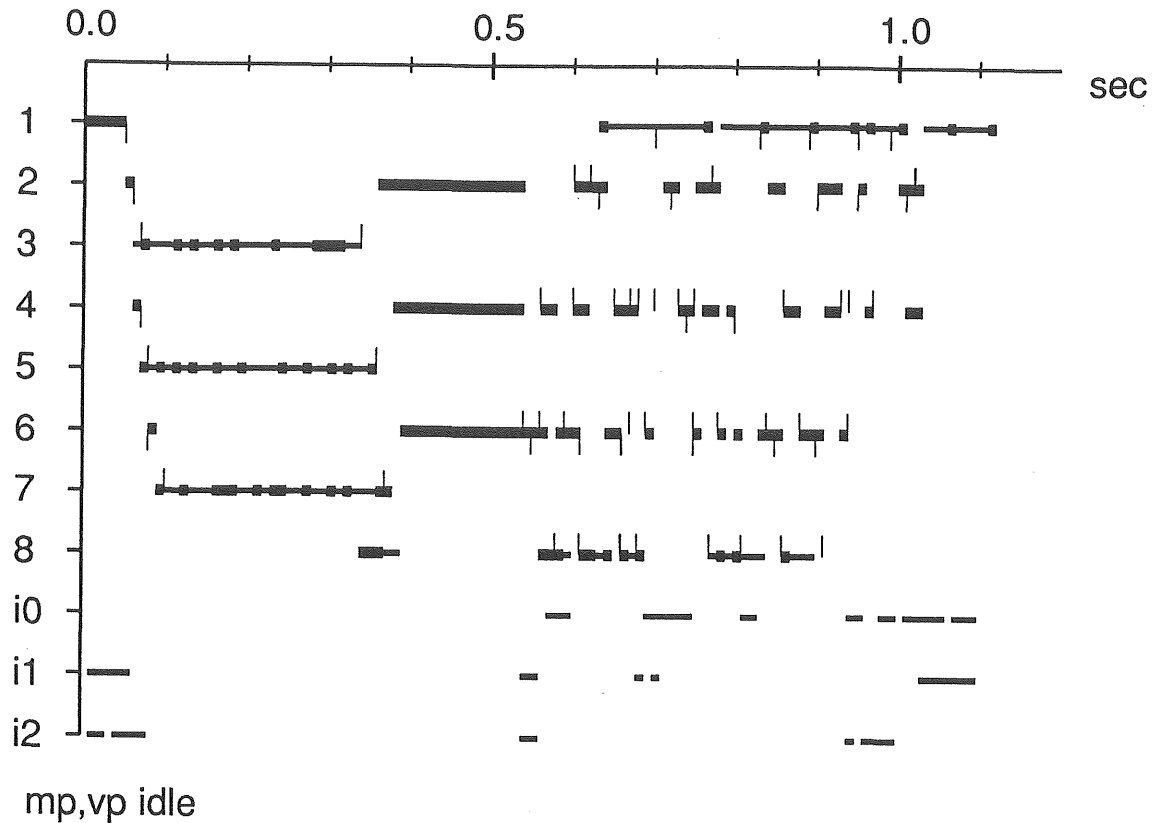


図3-6 実行の軌跡の記録を視覚化したもの

Figure 3-6: Visualized event traces of microprocess execution.

最初の4つは、それぞれ第2層の `mp_create()`、`mp_exit()`、`mp_sleep()`、`mp_wakeup()`、`mp_yield()` に関係している。ただし、`mp_wakeup()` において、既に実行可能になっているマイクロプロセスを再び実行可能にすることを試みた場合は、記録されない。(6)は、スケジューラにおいて、実行可能のマイクロプロセスに対してプロセッサが割り当てられたことを示すイベントである。(7)、(8)は、マイクロプロセスのCPU処理と入出力処理を区別して記録する必要がある場合に用いられる。(9)、(10)は、仮想プロセッサの動作、(11)は、プロセス全体の動作の示すイベントである。(12)、(13)は、軌跡を記録するためのファイル出力を表すイベントである。このようなイベントは、仮想プロセッサごとに記録される。

### 探り針効果

実行の軌跡を記録する場合、注意すべきことは、探り針効果 (probe effect) である。探り針効果とは、イベントを記録するために、その記録という行為がアプリケーションの実行に影響を与えてしまうことである。このマイクロプロセス・ライブラリでは、特に、イベントが発生した時刻を得るためのカーネル・コールのオーバーヘッドが大きい。これを改善するために、カーネル・コールによらない時刻の取得方法が求められる。

### 3. 6 仮想プロセッサを実現するカーネルの機能

仮想プロセッサは、カーネルにおいて実現される。利用者は、カーネル・コールを通じて、仮想プロセッサの機能を利用する。

#### 3. 6. 1 カーネル・コール

仮想プロセッサに関連したカーネル・コールには、次の4つがある。

##### (1) `vpid=vp_allocate(n)`

これは、新たに  $n - 1$  個の仮想プロセッサを生成することを要求するカーネル・コールである。このカーネル・コールの発行後、合計  $n$  個の仮想プロセッサによる並列処理が開始される。返り値は、仮想プロセッサの識別子である。

##### (2) `vp_switch(vpid)`

これは、同一プロセス内の他の仮想プロセッサへ制御を切り替えるためのカーネル・コールである。スピンロックを保持したまま実プロセッサを奪われた仮想プロセッサに制御を移すために用いられる。vpidは、次に実行すべき仮想プロセッサのヒントである。

##### (3) `vp_sleep(t)`

これは、その仮想プロセッサに対応している、一時的に不用になった実プロセッサをカーネルに返すためのカーネル・コールである。tは、実プロセッサを放棄する時間のヒントである。

##### (4) `vp_wakeup( n )`

これは、`vp_switch()`、あるいは、`vp_sleep()` で実行を停止している仮想プロセッサの実行再開を要求するカーネル・コールである。ただし、このカーネル・コールを用いなかったとしても、`vp_sleep()` を発行した仮想プロセッサは、指定された時間が経過した後は、実行を再開する。

仮想プロセッサを削除するカーネル・コールは、存在しない。それは、仮想プロセッサの数は、実プロセッサの数と同様に、プロセスの実行の途中で変化しないからである。

#### 3. 6. 2 カーネル・コールの利用例

図3-7、図3-8、図3-9に、C言語による上記のカーネル・コールの利用例を示す。(図3-7の1行目、2行目のキーワード `private`、図3-8の第1行目のキーワード `shared` については、次項(3.6.3項)において説明する。)

```

1: private int vpid ;      /* Virtual Processor ID */
2: private struct mpcb *mp_current ; /* Microprocess Control Block */
3:
4: main()
5: {
6:     init_single(); /* initialize in a sigle VP. */
7:     vpid = vp_allocate( 2 ); /* create a new VP. */
8:     init_multi(); /* initialize in multiple VPs. */
9:     mp_scheduler(); /* jump into mp_scheduler(). */
10:    /* Not reached. */
11: }

```

図3-7 vp\_allocate() カーネル・コールによる仮想プロセッサの生成  
Figure 3-7: Creatiting a virtual processor by invoking the vp\_allocate()  
kernel call.

図3-7第4行の main() は、このプログラムの実行が開始される場所を示している。プロセスには、起動時に1個の仮想プロセッサが自動的に割り当てられる。この仮想プロセッサにより、main() が実行される。まず、単一仮想プロセッサにより、さまざまな初期化を行う(第6行)。図3-7第7行では、vp\_allocate() カーネル・コールにより、新たに仮想プロセッサが1個作られ、これ以後、合計2個の仮想プロセッサが活動することになる。その後、複数の仮想プロセッサによる初期化を行い、制御をマイクロプロセス・スケジューラ(mp\_scheduler())に移す。その後、制御が再び戻って来ることはない。

図3-8は、簡単なマイクロプロセス・スケジューラの例である。この例では、第8行において、実行可能なマイクロプロセスを1個取り出し、それを mp\_run() により実行している。この時、実行可能なマイクロプロセスが存在しない場合、next\_mp が0になる。第11行において、vp\_sleep() カーネル・コールにより、現在この仮想プロセッサに対応している実プロセッサを開放している。VP\_SLEEP\_TIME は、実行を中断する時間のヒントである。

図3-9は、簡単なスピンロック(spin lock)の実現の例である。この例では、SPIN\_THRESHOLD (第5行)に定義された回数だけ test\_and\_set() (第8行)を繰り返し実行する。ここで、test\_and\_set() は、test-and-set 命令のような、原始にテストとセットを行う機械語命令で実現されるものとする。SPIN\_THRESHOLD に定義された数だけ回繰り返してもロックが保持できない場合、ループを抜けて、第13行目で vp\_switch() カーネル・コールを発行している。カーネルは、プロセッサの横取りを行った仮想プロセッサが存在しないかを調べ、それが見つかった時には、その仮想プロセッサへ制御を移す。

```

1: shared lock_t mp_ready_lock ;
2:
3: mp_scheduler()
4: {
5:     struct mpcb *next_mp ;
6:
7: loop:  lock( &mp_ready_lock );
8:         next_mp = mp_ready_dequeue();
9:         unlock( &mp_ready_lock );
10:
11:        if( next_mp == 0 )
12:        {
13:            vp_sleep( VP_SLEEP_TIME );
14:            goto loop;
15:        }
16:        mp_run( next_mp );
17:        /* Not reached. */
18: }

```

図 3 - 8 vp\_sleep() カーネル・コールによる実プロセッサの開放

Figure 3-8: Releasing a real processor by invoking the vp\_sleep()  
kernel call.

```

1: lock( lock_t *v )
2: {
3:     int count ;
4: loop:
5:     count = SPIN_THRESHOLD ;
6:     do
7:     {
8:         if( test_and_set( v ) == 0 )
9:             return;
10:    }
11:    while( --count > 0 );
12:
13:    vp_switch( VP_NONE );
14:    goto loop;
15: }
16:
17: unlock( lock_t *v )
18: {
19:     *v = 0 ;
20: }

```

図 3 - 9 vp\_switch() カーネル・コールによる他の仮想プロセッサへの制御の移動

Figure 3-9: Forcing context switch to another virtual processor by  
invoking the vp\_switch() kernel call.

### 3. 6. 3 仮想プロセッサの固有領域

本方式では、仮想プロセッサごとに固有のメモリ領域がある。すなわち、各仮想プロセッサには、アドレス空間の一部に、仮想プロセッサごとに異なる物理メモリが参照される領域が存在する。これを仮想プロセッサの固有領域 (private area) とよぶ。これに対して、固有領域以外の領域を共有領域 (shared area) とよぶことにする。固有領域は、マイクロプロセスを実現するライブラリにおいて、仮想プロセッサごとの情報を格納するために利用される。この領域は、実プロセッサのレジスタが拡張されたものと考えることができる。

図3-10に、1つのプロセスに属する2つの仮想プロセッサの論理アドレス空間と物理メモリの対応を示す。これは、図3-7に示したプログラムの実行の結果、構築される対応である。図3-7第7行の `vp_allocate()` カーネル・コールにより、新たに仮想プロセッサが1個作られ、合計2つの仮想プロセッサが活動している。

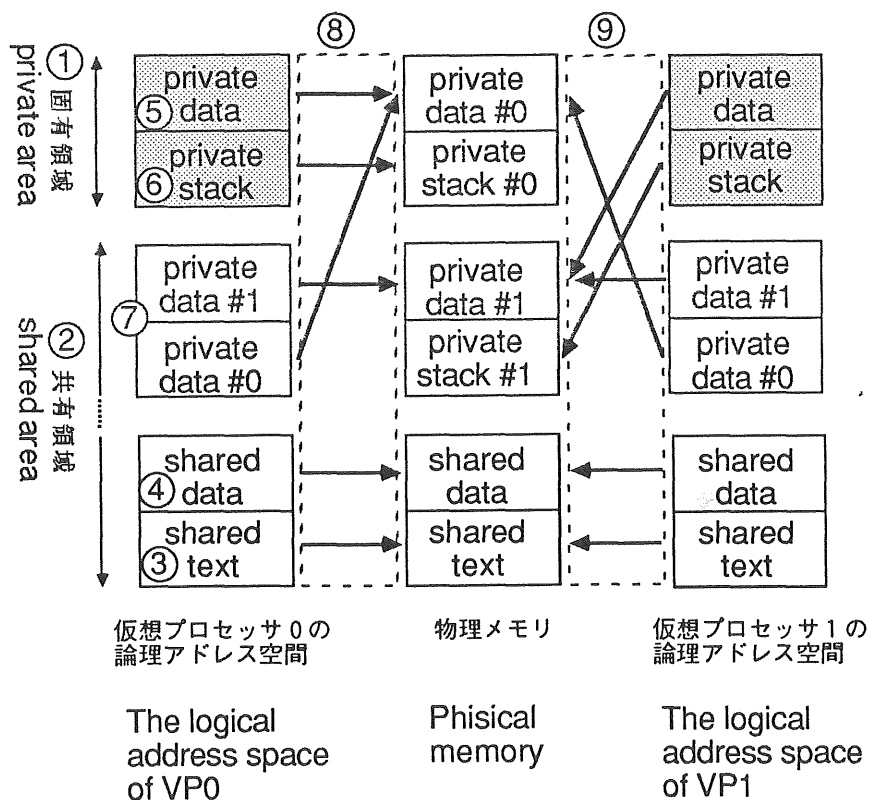


図3-10 2つの仮想プロセッサの論理アドレス空間と物理メモリの対応  
Figure 3-10: Mapping the logical address spaces of two virtual processors to physical memory.

図3-10では、アドレス空間の上部 (アドレスが大きい方、図でも上) に固有領域 (①)、アドレス空間の下部に共有領域 (②) が割り付けられている。共有領域は、



テキスト・セグメント (text segment) (③)、データ・セグメント (data segment) (④) に分けられる。前者は、機械語命令を格納する部分であり、読み専用になっている。後者は、データを格納する部分であり、読み書き可能になっている。固有領域は、固有データ・セグメント (private data segment) (⑤) と固有スタック・セグメント (private stack segment) (⑥) に分けられ、両方とも読み書き可能である。前者は、仮想プロセッサごとに固有の変数を割り付けるために利用される。後者は、マイクロプロセス・スケジューラ (図 3-7 第 9 行の mp\_scheduler()) のスタックとして用いられる。各仮想プロセッサの固有データ・セグメントは、他の仮想プロセッサによる参照を可能にするために、アドレス空間の一部に割り付けられている (⑦)。

図 3-7 において、キーワード private (第 1 行、第 2 行) は、この変数を固有領域に割り付けるように指示するものである。これとは逆に、変数を共有領域に割り付けるように指示するキーワードは、図 3-8 第 1 行に現れている shared である。このような機能は、D y n i x システム [6] におけるいくつかの言語処理系において既に実現されている。D y n i x では、U N I X のオブジェクト・コード形式 (a.out 形式) を拡張して、共有と固有の 2 つのデータ・セグメントを提供している。

### 3. 6. 4 仮想プロセッサの固有領域の利用

仮想プロセッサの固有領域は、次の 4 種類のデータを格納するために用いられる。下記 (1)、(2)、(3) は、固有データ・セグメントに、(4) は、固有スタックセグメントに割り付けられる。

(1) 仮想プロセッサの識別子 (図 3-7 第 1 行)。

(2) 仮想プロセッサごとに必要となるデータ。

利用者空間において軽量プロセスを実現する場合、このようなデータが必ず必要になる。たとえば、カーネル・コールのエラーを保持する変数 `errno` や現在実行しているマイクロプロセスの制御ブロック (`mpcb`、microprocess control block) へのポインタが仮想プロセッサごとに必要となる (図 3-7 第 2 行)。

(3) 共有する必要がない静的変数。

C 言語のライブラリ関数の中には、静的変数 (static variable) に結果を代入しそのポインタを返すものや、静的変数に中間状態を格納するものが多数存在する。前者の例として、時刻を文字列に変換する `ctime()`、各種検索用ライブラリ `gethostbyname()`、`getpwent()`、後者の例としては、引数を解析するルーチン `getopt()` があげられる。このような静的変数を固有領域に割り付けるにより、逐次プログラム用に開発されたライブラリ関数を並列プログラムにおいて利用することが可能となる。

(4) マイクロプロセス・スケジューラ用のスタック。

これは、(2)の特殊なものである。(どの仮想プロセッサも任意のマイクロプロセスを実行することを可能にするために、マイクロプロセスのスタックは、共有領域に置かれる。)

固有スタック・セグメントの存在により、図3-7に示したようなプログラムをスタック・ポインタの切替えなしに動作させることが可能になっている。たとえば、図3-7の第7行において、カーネル・コール `vp_allocate()` を発行する仮想プロセッサ(仮想プロセッサ0)も、新たに生成される仮想プロセッサ(仮想プロセッサ1)も、同一のアドレスに復帰する。そして、どの仮想プロセッサも同一のアドレスにあるスタックを利用することが可能となる。

### 3.6.5 固有領域を利用しない方法との比較

仮想プロセッサごとの固有領域を利用しない方法として、仮想プロセッサの識別子を実プロセッサのレジスタに格納する方法が考えられる。この方法では、3.6.4項で述べた(2)、(3)、(4)を、共有領域にある配列とて実現しなければならない。たとえば、現在実行中のマイクロプロセス制御ブロック(struct `mpcb`)へのポインタ `mp_current` は、図3-11に示すように共有領域にある配列を用いて実現することができる。ここで、`MAX_NVP`(図3-11第1行)は、仮想プロセッサ数の最大値である。`get_vpid()`(図3-11第3行)は、実プロセッサに格納された仮想プロセッサの識別子を参照するための関数である。このプログラムでは、マクロを用いて配列を単純な変数として扱えるようにしている。

```
1: #define MAX_NVP 32
2: struct mpcb *mp_current_array[MAX_NVP]; /* shared */
3: #define mp_current mp_current_array[get_vpid()]
4: ....
5: {
6:     mp_lock( mp_current );
7:     ....
8:     mp_unlock( mp_current );
9: }
```

図3-11 配列と仮想プロセッサ識別子による  
仮想プロセッサごとの固有領域の実現

Figure 3-11: A per-virtual-processor variable by using  
virtual processor identifiers and an array

配列と実プロセッサのレジスタに格納された仮想プロセッサの識別子を用いる方法には、次のような問題点がある。

- (1) 配列のアドレス計算のオーバーヘッドにより実行速度が低下する。
- (2) C言語のように静的に配列の大きさを決定する言語では、最大の仮想プロセッサ数分のメモリが常に必要になり、メモリの無駄が生じる。メモリの無駄を避ける方法として、必要なメモリをヒープ上に動的に確保する方法が考えられる（C言語では、`malloc()` ライブラリ関数を用いる）。しかしながら、この方法では、メモリへのアクセス回数が増加し、実行速度の低下を招く。
- (3) 実プロセッサのレジスタを1個専有することになるため、応用プログラムの利用可能レジスタ数が減る。
- (4) 大域的なレジスタ変数機能を提供している言語処理系が一般的に利用可能にはなっていない。

仮想プロセッサの固有領域を用いることにより、これらの問題を解決することができる。

仮想プロセッサの識別子を実プロセッサのレジスタに格納する代りに、必要になった時にカーネル・コールを発行する方法も考えられる。しかしながら、利用者空間内において軽量プロセスを実現しているという利点が失われてしまうので、この方法を採用することはできない。

### 3.7 マイクロプロセスと仮想プロセッサの利用

この節では、マイクロプロセスと仮想プロセッサの利用法を示す。最も重要な利用は、共有メモリ型マルチプロセッサにおける並列処理である。そのほかに、ファイルの先読み、入出力処理とCPU処理の重ね合わせなどにも利用される。

#### 3.7.1 共有メモリ型マルチプロセッサにおける並列処理

ここでは、簡単な例として、ループ分割（loop splitting）をマイクロプロセスを利用して実現する。図3-12に示した、配列の要素の和を求める逐次プログラムを、マイクロプロセスを用いてループ分割したものを、図3-13に示す。ここでは、同期プリミティブとして計数セマフォ（counting semaphore）を用いている。

図3-13で、`nmprocs` は、マイクロプロセスの数である。CPU処理を主とする（CPU-intensive）プログラムの場合、実プロセッサ数と等しい数の仮想プロセッサを生成する。そして、実プロセッサ数の数倍のマイクロプロセスを生成して、並列処理を行う。ここで、実プロセッサ数よりも多くのマイクロプロセスを生成する点に、多重プログラミング・システムの特徴が現れている。

```

1: array_sum( a, n )
2:   int a[];
3:   int n ;
4: {
5:   int i, sum ;
6:     sum = 0 ;
7:   for( i=0 ; i<n ; i++ )
8:     {
9:       sum += a[i] ;
10:    }
11:   return( sum );
12: }

```

図 3 - 1 2 ループ分割する前の逐次プログラム

Figure 3-13: A sequential program to be loop-split.

そのサイトの負荷が低い場合、このプログラムは、単一プログラミング・システムとほぼ同等の速度で実行される。そのサイトの負荷が大きく、割り当てられる実プロセッサの数が少なくなったとしても、このプログラムは、動作する。さらに、並列処理の途中で、割り当てられる実プロセッサの数が変動したとしても、まったく問題なく動作する。

単一プログラミング・システム用の並列プログラムの場合、プロセッサ数に等しいプロセス（並列処理の単位）が生成される。すなわち、`nmprocs` を実プロセッサ数に合わせる。このようなプログラムをそのまま多重プログラミング・システムに移植しても、効率よく動作しない。それは、全ての実プロセッサがそのプログラムに割り当てられるとは限らないからである。上のプログラムにおいて、`nmprocs` を実プロセッサ数にした時、プロセッサが1個欠け、`nmprocs-1` 個しか割り当てられなかったとすると、そのプログラムの実行時間は、`nmprocs` 個のプロセッサが割り当てられた時の2倍になる。さらに、同期処理において、ビジー・ウェイトが用いられている場合、大量のCPU時間が浪費されてしまう。この浪費は、多重プログラミング・システムでは許されない。

### 3. 7. 2 入出力の重ね合わせ

複数の入出力の並列実行や、ディスク・ブロックの先読みを実現するためにも、マイクロプロセスと仮想プロセッサを利用することができる。たとえUNIXでは、ファイルをシーケンシャルにアクセスする場合に限り、1ブロックの先読みを行う。データベース処理では、ファイルをランダムにアクセスする場合でも、次に必要となるディスク・ブロックを予め知ることができる場合が非常に多い。このような場合、以下に述べるような方法により、先読みを実現することができる。

```

1: array_sum( a, n, nmprocs )
2:   int a[];   /* int a[0..n-1]; */
3:   int n ;
4:   int nmprocs ; /* number of microprocesses */
5: {
6:   int i, sum, bite ;
7:   sem_t mutex ;
8:   sem_t done ; /* private semaphore */
9:
10:   sem_init( &done, -nmprocs+1 );
11:   sem_init( &mutex, 1 );
12:   bite = n / nmprocs ;
13:   for( i=0 ; i<n ; i+=bite )
14:   {
15:       mp_create( array_sum_task, ARRAY_SUM_TASK_SSIZE, MP_READY,
16:                 &a[i], bite, &mutex, &done, &sum );
17:   }
18:   sem_P( &done );
19:   return( sum );
20: }
21:
22: array_sum_task( a, n, mutexp, donep, sump )
23:   int a[];
24:   int n ;
25:   sem_t *mutexp ;
26:   sem_t *donep ;
27:   int *sump ;
28: {
29:   int i, sum ;
30:   sum = 0 ;
31:   for( i=0 ; i<n ; i++ )
32:   {
33:       sum += a[i] ;
34:   }
35:   sem_P( mutexp );
36:   *sump += sum ;
37:   sem_V( mutexp ); /* mutual exclusion */
38:   sem_V( donep ); /* wakeup parent */
39:   mp_exit();
40: }

```

図 3 - 1 3 マイクロプロセスによるループ分割

Figure 3-13: Loop-splitting by using microprocesses.

図 3 - 1 4 に示すように、データ処理を行うマイクロプロセスの他に、データ入力を行うマイクロプロセスとバッファを 2 個用意する。データ入力を行うマイクロプロセスは、データ処理を行うマイクロプロセスが一方のバッファについて処理を行っている間に、別のバッファに次に必要となるブロックの読み込みを行う。この時、データ入力を行うマイクロプロセスを実行していた仮想プロセッサは、入出力の完了待ち状態になる。しかしながら、そのプロセスには、他に実行可能な仮想プロセッサが存在

するので、その仮想プロセッサにコンテキストを切り替えることで、即座に応用プログラムの実行に戻ることが可能である。

従来、ファイルの先読みは、非同期入出力機能を利用して実現されることが多かった。しかしながら、非同期入出力は、同期型入出力に比較してプログラミングが困難であるという問題がある。マイクロプロセス（軽量プロセス）を利用する場合、論理的に並列に実行したい部分をマイクロプロセスとして記述するだけで、自然に並列性が抽出される。

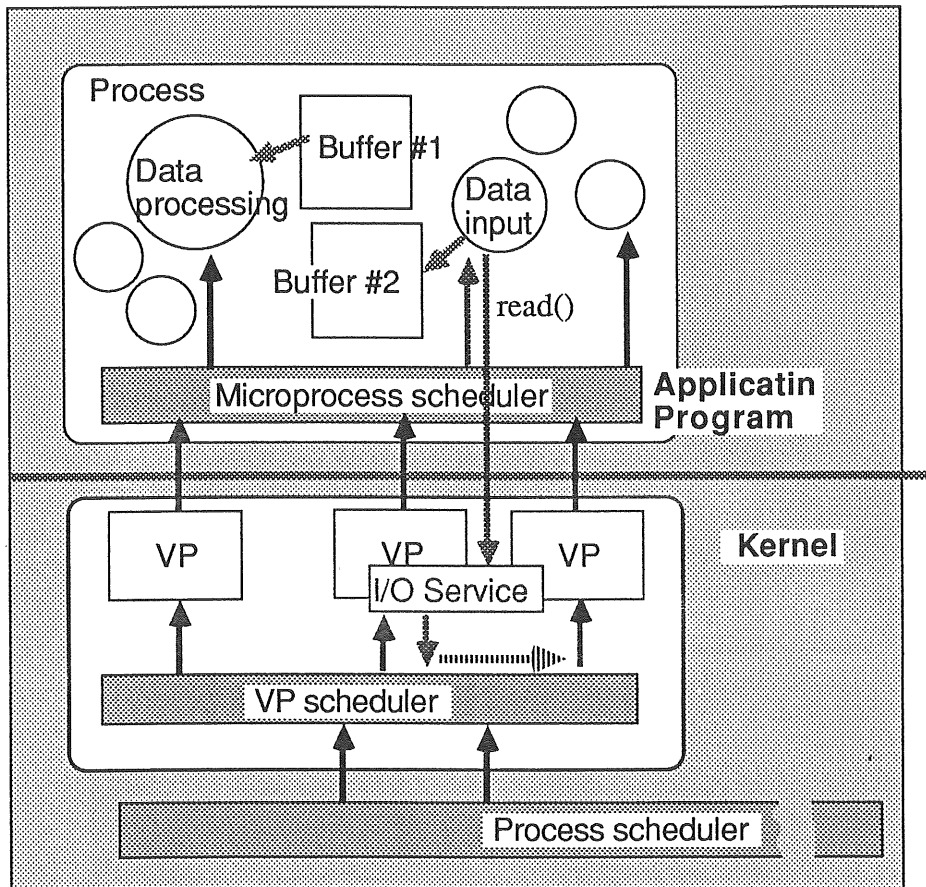


図3-14 マイクロプロセスによる入出力の重ね合わせ

Figure 3-14: Overlapping I/O and CPU processing by using microprocesses.

### 3.7.3 非同期通信の実現

本システムでは、プロセス間通信のプリミティブとして、同期式のもの（遠隔手続き呼出し）を採用している。応用プログラムにおいて、非同期式のものが必要になった場合は、マイクロプロセスと仮想プロセッサを用いて簡単に実現することができる。

非同期処理を行いたい部分で、新たにマイクロプロセスを生成する。新たに生成されたマイクロプロセスは、他の仮想プロセッサにより実行される。その後、必要になった時点で、先ほどのマイクロプロセスの実行完了を待つ。

### 3.7.4 動的負荷分散

分散型オペレーティング・システムの研究分野において、動的負荷分散は、重要な課題となっている。その手段の1つとしてプロセスの移動 (process migration) があげられる。プロセスの移動とは、実行中のプロセスを別のサイトに移動し、プロセスの実行を継続することである (図3-15(a))。

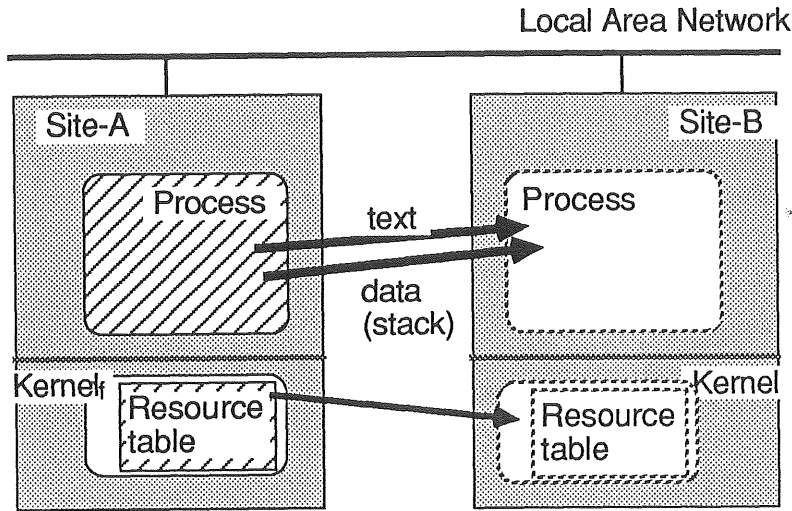


図3-15(a) 重量プロセスの移動

Figure 3-15(a): Migration of a heavyweight process.

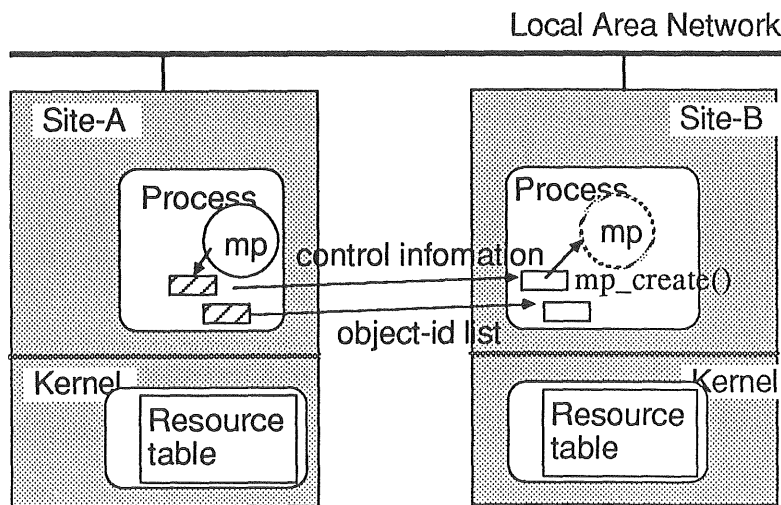


図3-15(b) 軽量プロセスの移動

Figure 3-15(b): Migration of a lightweight process.

図3-15 プロセスの移動

Figure 3-15: Process migration.

プロセスの移動において、重量プロセス（資源割当てと保護の単位）ではなく、軽量プロセスを用いる方法が考えられる（図3-15（b））。重量プロセスの場合、プロセスの移動そのものに要するコストが大きく、かつ、時間もかかる。これに対して、軽量プロセスの場合、移動に必要となるデータ量を削減することが可能で、処理効率が改善されると思われる。その理由は、第1に、複数のサイトを利用して1つの応用プログラムを実行する場合、各サイトで実行されるプログラムのロード・モジュールが同一であることが多いからである。この場合、移動先のプロセスに同一のテキストが存在するので、テキストを転送する必要がない。

第2に、マイクロプロセスの性質によっては、非常に小さな制御情報からマイクロプロセスを再構成することが可能な場合がある。この場合、その制御情報だけを送ることにより、マイクロプロセスの固有データ領域全体を転送する必要がなくなる。このように、応用プログラムの性質を利用することにより、応用プログラムのレベルにおいてマイクロプロセスの移動を行い、動的負荷分散を実現することも可能である。

小さな制御情報からマイクロプロセスを再生することは、ARGUSにおけるガーディアン（guardian）の再生と類似している[49]。ガーディアンとは、一種のプロセスであり、内部に変数を持ち、ハンドラ（handler）と呼ばれる手続きを通じてのみ内部の変数へのアクセスが許される。内部の変数は、通常のメモリ中のみ記録され、障害が発生すると失われる揮発性の変数と、ディスクなどに保存され、障害が発生しても失われない安定な（stable）変数に分類される。障害が発生すると、ARGUSでは、この小さな安定な変数から、通常の変数の状態を回復し、ガーディアンを再生する。このような機構は、マイクロプロセスの移動の目的も活用することができる。

応用プログラムのレベルにおいてマイクロプロセスの移動を行う場合、カーネル中のプロセスの資源管理表の移動がしばしば問題となる。本システムでは、この問題は生じない。それは、結合を作らない（connection-less）の通信プリミティブである遠隔手続き呼出しと、利用者レベルで受け渡し可能なオブジェクト識別子を利用しているからである。現在参照しているオブジェクトのオブジェクト識別子を移動先に転送することにより、プロセスの資源表を転送せずに元のプロセスと同一のオブジェクトを参照することが可能である。

### 3. 8 応用固有の同期・通信プリミティブの開発

#### 3. 8. 1 層構造を持つライブラリの利用

3. 5節で述べたように、本システムのマイクロプロセス・ライブラリ（利用者レベルの軽量プロセスを実現するライブラリ）は、層構造を持っている。このライブラリを利用する場合、各並列プログラムのプログラマは、応用固有の同期・通信プリミティブを構築するために、たとえば、以下で示すようなレベルで機能を利用すること



ができる。高いレベルの機能を利用するほど、細かい記述から開放される。一方、低いレベルの機能を利用するほど、高い性能を得ることができる。

- (レベル4) 第3層の同期・通信プリミティブから選択して利用する。この場合、複数のプリミティブを混在させて利用することも可能である。
- (レベル3) 第3層の同期・通信プリミティブを利用して、独自の同期・通信プリミティブを開発する。
- (レベル2) 第2層の sleep/wakeup モデルに基づく同期・通信プリミティブを利用して、独自のプリミティブを開発する。
- (レベル1) 第1層のスピンロック・ライブラリとコンテキストの保存・回復の機能を用いて、独自のプリミティブを開発する。このレベルでは、アーキテクチャから独立したプログラムを開発することができる。
- (レベル0) 第0層のカーネル・コールだけを利用して同期・通信プリミティブを開発する。

### 3. 8. 2 コルーチンの性質の利用

マイクロプロセス間の同期・通信プリミティブは、応用プログラムごとに開発される。本方式では、コルーチンの性質を活用して、固有の同期・通信プリミティブを効率的に開発することができる。

3. 4. 1 項、および、3. 5. 2 項で述べたように、マイクロプロセスは、コルーチンと非常に類似している。相違点は、マイクロプロセスの場合、並列処理の対象となる点にある。そのためマイクロプロセスの場合、マイクロプロセスの共有変数の相互排除を行う必要がある。逆に、マイクロプロセスにおいて、共有変数の相互排除操作を削除することにより、完全にコルーチンに還元することが可能である。

並列プログラムのデバッグにおいて、実行の再現性がないことが大きな障害となっている。対象的に、コルーチンの場合、プログラムの実行に再現性がある。これは、プログラムが逐次的に実行されるからである。この再現性があるという性質は、プログラムのデバッグにおいて、非常に都合がよい。本システムでは、1 個の仮想プロセッサを用いることにより、容易にコルーチンと同じ動作を行わせることができる。すなわち、カーネル・コール `vp_allocate()` の引数として 1 を指定するか、または、そのカーネル・コールを発行しなければよい。

以上のことより、本マイクロプロセス方式において、新たにマイクロプロセス間の同期・通信プリミティブを開発する場合、次の様な手順で効率的に行うことができる。

- (1) 相互排除操作を含まないマイクロプロセスにより、同期・通信プリミティブを構築する。これは、コルーチンと完全に等価である。
- (2) マイクロプロセス間の共有変数に対する相互排除の操作を付加する。そして

これを、1つの仮想プロセッサを用いて実行する。この時に、スピン・ロックの解除を忘れる、あるいは、1つのマイクロプロセスで同じ変数に2回ロックを試みるなどのバグが取除かれる。この時点で、実行のタイミングに関するバグを除いて、ほとんどのバグが取除かれる。

(3) 最後に、複数の仮想プロセッサを用いて実行し、タイミングに関するバグを取る。

ステップ(2)において、スピン・ロックを実現するライブラリ手続きとして、デバッグ用の特殊なものを用意している。それは、ロック操作において、既にロックされている場合に、メッセージを表示して直ちに処理を中断するものである。1つの仮想プロセッサにおいては、ロックが重なることは、デッドロックを意味する。

### 3. 8. 3 応用固有の同期・通信プリミティブの開発例

この節では、データベースの並列処理システムとして当研究室で開発しているSMASHを取上げ、本システムで提案した軽量プロセス機能の有効性について論じる。SMASHは、データベース、および、知識ベースを対象とした並列処理システムである[38][39][40][50]。SMASHが対象とするハードウェアは、本システムと同様、図1-1に示すような共有メモリ型マルチプロセッサと単一プロセッサが高速ネットワークにより結合された環境である。

SMASHの特徴は、関数型プログラミングの概念をデータベース処理に適応している点にある。関数の評価方式としては、要求駆動型評価(demand driven evaluation)を用いている。SMASHでは、関数間のデータの受け渡しにストリームを用いている。SMASHでは、任意のデータベース演算を関数として記述し、システム内に組み込むことを可能にしている。このため、関数間に存在する並列性を抽出することを可能になり、データベースの多様な応用分野に柔軟に対応することが可能となっている[38][57][73]。

#### 3. 8. 3. 1 SMASHのプリミティブ・セット

SMASHは、1章において述べた単一プログラミング・システムに相当する。SMASHは、図1-1のような並列/分散処理ハードウェア環境を、専有して使用する。

SMASHでは、このような環境のもとで、関数型計算を実行するためのプリミティブ・セットを設定している。これは、抽象計算機の命令セットであり、一種の中間言語である。そして、各プリミティブに対応する実行時ルーチンが存在する。関数型言語で記述されたプログラムは、このプリミティブを呼び出すコードを含むオブジェクト・コードにコンパイルされる。したがって、このプリミティブ・セットを実現することが、SMASHシステムを移植することになる。ここでは、SMASHのプリ

ミティブ・セットを本システムが提供する機能を利用して実現することを考える。

S M A S Hのプリミティブの主要な機能は、次の3つである。

- (1) 関数のインスタンスの生成・消去
- (2) ストリーム・データの通信路であるチャンネルの生成・消去
- (3) チャンネルを通じたストリーム・データの受け渡し

S M A S Hにおけるフロー制御 (flow control) は、要求駆動型評価に基づいておこなわれる。すなわち、ストリーム・データの生産者側の実行は、消費者側からの要求 (demand) により制御される。これは、T C P / I Pにおけるウインドウ制御と類似している [46] [101]。S M A S Hでは、高いレベルにおいてフロー制御を行っている点に特徴がある。

主なプリミティブの名称と機能を、以下に示す。

`new( func, site, arg1, arg2, ... )`

指定された関数 `func` のインスタンス (これを関数インスタンスとよぶ) をサイト `site` 上に生成する。`arg1, arg2, ...` は、関数インスタンスへ引数として渡される。引数として、以下で述べるチャンネルの識別子を渡すことができる。この時、自動的にコネクションが形成される。

`channel_t channel( type, gran , prod_site, cons_sites )`

チャンネルと呼ばれるストリーム通信路を設定する。`type` は、ストリーム要素の型、`gran` は、チャンネル用のバッファの大きさを指定する。`prod_site` は、ストリームの生産者となる関数インスタンスが存在するサイト、`cons_sites` は、消費者となる関数インスタンスが存在するサイトを指定する。消費者として、複数のサイトを記述することができる。

`get( p, data )`

`receive( p, data )`

`put( p, data )`

`send( p, data )`

`mark_eos( p )`

`get()`、`receve()` は、ストリームからデータを入力するプリミティブ、`put()`、`send()` は、ストリームへデータを出力するプリミティブである。`p` は、ポートと呼ばれるコネクションが形成されたチャンネルである。`get()` と `put()` は、バッファが、あるいは、データが未到着で入出力に失敗した際、そのまま待ち続ける。これに対して、`receive()` と `send()` は、失敗を返し、ブロックされることはない。`mark_eos()` は、ストリームの終了の目印をつける。

`select( p1, p2, ... )`

`enable( p )`

`disable( p )`

`select` は、入出力可能なストリームの非決定的な選択を行う。UNIX 4.3 BSD の `select` システム・コールと同様な機能を提供する [46]。

`enable()/disable()` は、指定されたポートを選択の対象に入れる／対象から外す操作を行う。

### 3. 8. 3. 2 1 サイト内における SMASH のプリミティブの実現

前項で述べたプリミティブ・セットを、本システムが提供する軽量プロセス機能を利用して実現する方法について述べる。ここでは、共有メモリ型マルチプロセッサにおける実現方式を述べ、ネットワーク環境における実現は次の 3. 8. 2. 3 で述べる。ここで、単一プロセッサは、プロセッサ数が 1 の共有メモリ型マルチプロセッサと考える。

基本的な実現方式は、SMASH の関数インスタンスをマイクロプロセスとして実現することである。関数インスタンスの生成 `new()` プリミティブは、マイクロプロセスを生成するライブラリ関数 `mp_create()` を用いて実現した。

関数インスタンス間の通信プリミティブを、本マイクロプロセス・ライブラリの第 2 層の機能を利用して第 3 層に実現した。SMASH のチャンネルは、図 3-16 に示すように、マイクロプロセス間の共有データ、マイクロプロセスの固有データ、実際にストリームデータを格納するバッファにより実現される。それぞれ、チャンネル制御表、ポート制御表、ページとよぶ。チャンネル制御表は、複数のマイクロプロセスから同時にアクセスされる可能性があるため、ロックにより相互排除が行われる。ポート制御表には、1 つのマイクロプロセスに固有の情報が格納される。したがって、相互排除を行う必要がない。1 つのページについて処理を行っている間、マイクロプロセスは、ポート制御表だけを参照し、チャンネル制御表を参照することはない。このようにして、マイクロプロセス間の同期制御のオーバーヘッドを軽減している。

図 3-16 では、ストリームの通信路であるチャンネル 1 個に対して、2 個のページが割り当てられている。この結果、ストリームの生産者となるマイクロプロセスと消費者となるマイクロプロセスを並列に実行することが可能となっている。

SMASH プリミティブの実装を、C 言語により行った。ソース・プログラムのお大きさは、全体で約 2500 行である。

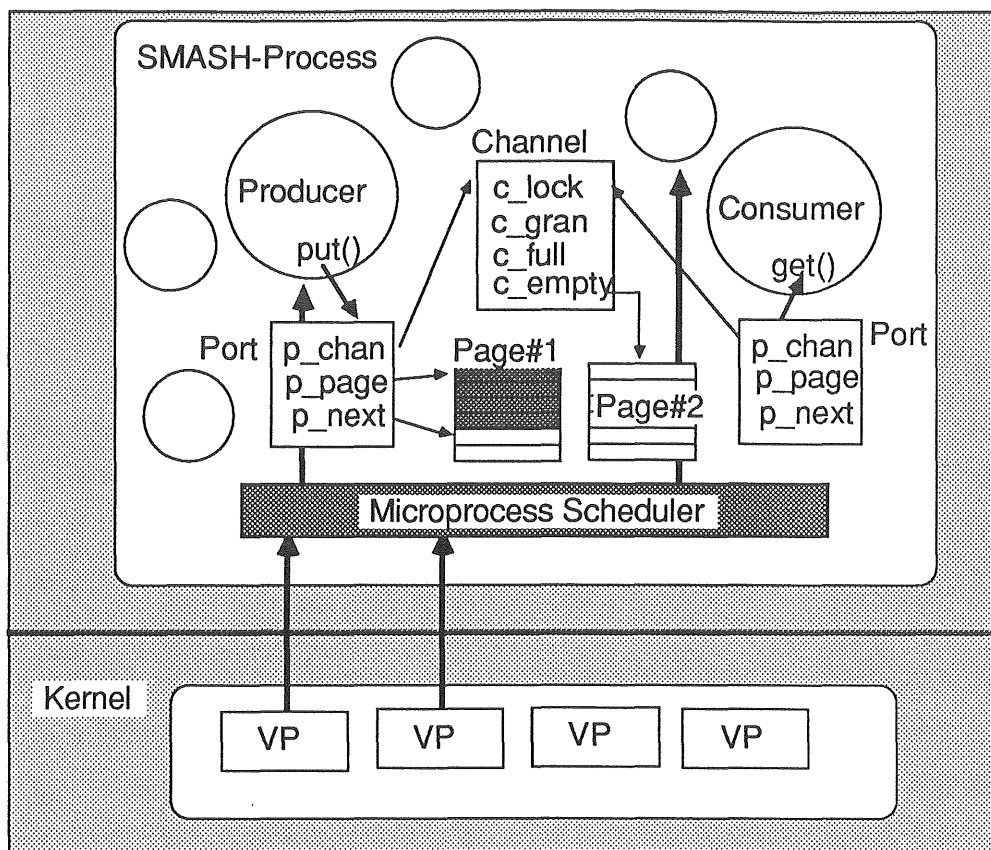


図 3 - 1 6 共有メモリ型マルチプロセッサにおける  
マイクロプロセスを利用した S M A S H の実現

Figure 3-16: The implementation of SMASH primitives on a shared-memory multiprocessor by using microprocesses.

### 3. 8. 3. 3 ネットワーク環境における S M A S H のプリミティブの実現

ネットワーク環境における S M A S H のプリミティブ・セットの実現を、図 3 - 1 7 に示す [ 7 2 ]。S M A S H のインスタンスは、1 サイト内と同様に、マイクロプロセスとして実現される。マイクロプロセスの生成に関しては、完全にアプリケーションの制御にまかされており、本システムのカーネルは、一切関与しない。したがって、遠隔サイト上に関数インスタンスを生成する場合、ネットワーク通信機能を用いて、遠隔サイト上の S M A S H プロセスに対してメッセージを送る。メッセージを受取ったプロセスは、自分自身の中にマイクロプロセスを生成する。

この実現においては、関数インスタンスの他に、ネットワーク通信を実現するためにマイクロプロセスが利用されている。具体的には、ネットワークからのメッセージを受信するために、マイクロプロセスが利用されている。マイクロプロセス（軽量プロセス）の利用により、並列性の記述が自然に行われ、プログラムの論理的な構造が整理されている。

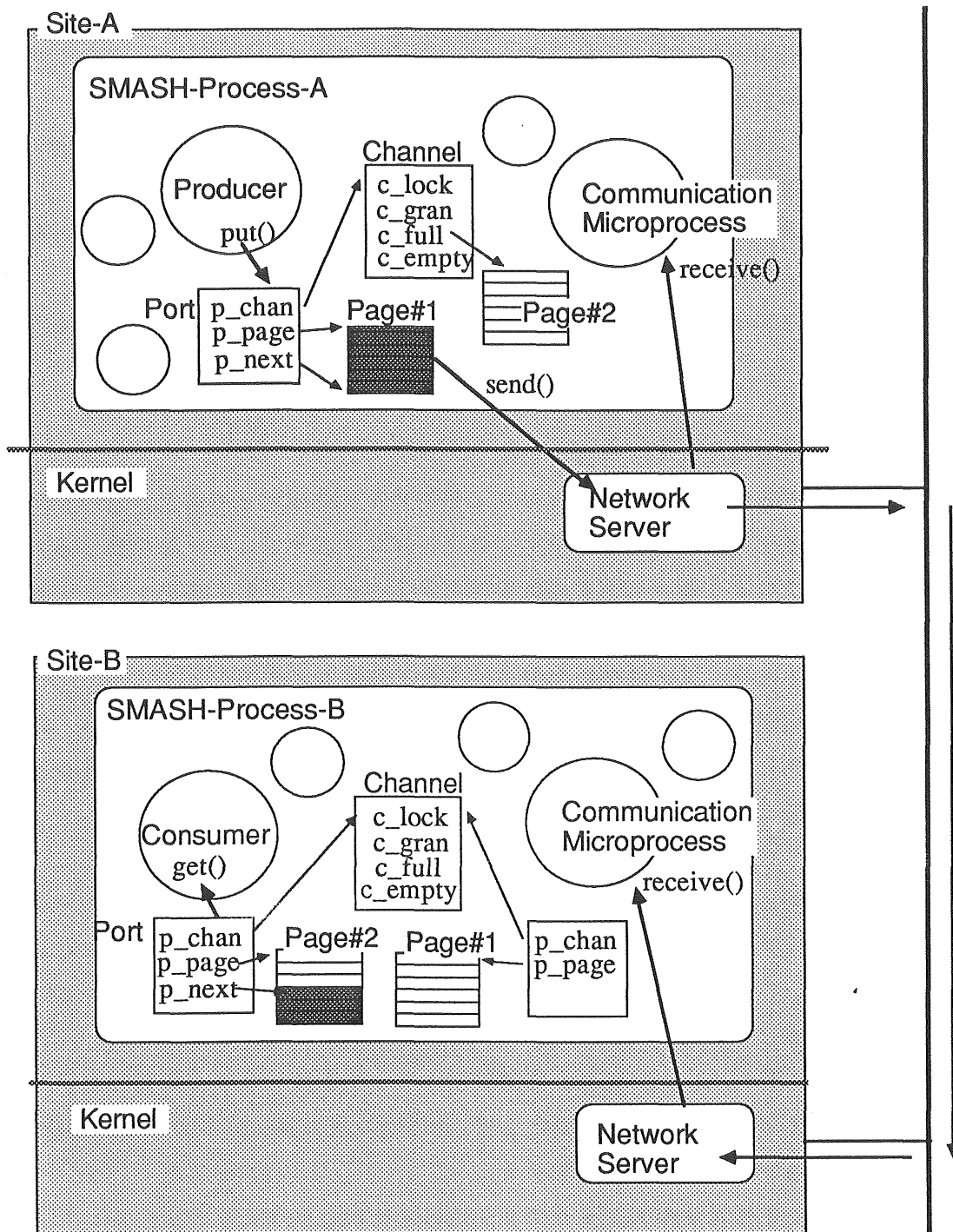


図 3 - 1 7 ネットワーク環境における  
マイクロプロセスを利用したSMASHの実現

Figure 3-17: The implementation of SMASH primitives on networked processors by using microprocessors.

### 3. 9 応用固有のスケジューラの開発

本方式では、応用固有のスケジューラを、マイクロプロセス・スケジューラとして、利用者空間内に構築する。この方法の利点は、以下の通りである。

- (1) 応用固有の情報の入手が容易である。スケジューラは、プロセス間の保護の壁を越えることなく、実行の状況を調べることができる。
- (2) カーネル中にスケジューラを設ける方法と比較して、スケジューラを取り替えることが容易である。
- (3) 同時に動作する複数の並列応用プログラムのスケジューラが互に干渉しあうことがない。

#### 3. 9. 1 層構造を持つライブラリの利用

3. 5 節において述べたマイクロプロセス・ライブラリを利用する場合、各並列プログラムのプログラマは、応用固有のマイクロプロセス・スケジューラを構築するために、たとえば、次のようなレベルで機能を利用することができる。高いレベルの機能を利用するほど、細かいスケジューリングの記述から開放される。一方、低いレベルの機能を利用するほど、きめ細かい制御を行うことができる。

- (レベル5) 予め用意されているスケジューラから選択して利用する。(現在、代表的なスケジューリング方式として、F I F O、および、固定優先順位方式を用意している。)
- (レベル4) マイクロプロセス・ライブラリ中のレディ・キュー・モジュール  
(図3-3)と同じ外部仕様をもつモジュールを記述し、残りのマイクロプロセス・ライブラリと結合して利用する。3. 9. 2 項において、このレベルを用いて記述した優先順位式スケジューラについて述べる。
- (レベル3) 第2層の sleep/wakeup 同期プリミティブと同じレベルに新たなプリミティブを追加する。たとえば、sleep\_and\_wakeup のように、複合的なプリミティブを追加する。
- (レベル2) マイクロプロセス・ライブラリ中のスケジューラ・モジュール  
(図3-3)と同じ外部仕様をもつモジュールを記述し、残りのマイクロプロセス・ライブラリと結合して利用する。
- (レベル1) 第1層を利用して、スケジューラを記述する。
- (レベル0) 第0層(カーネル・コール)だけを利用して、スケジューラを記述する。

#### 3. 9. 2 応用固有の軽量プロセス・スケジューラの開発例

ここでは、3. 8. 3 項において述べた S M A S H システムにおけるスケジューラの記述を例に、応用固有のスケジューラの開発について述べる。

次の様な3種類のスケジューラを記述してみた。

(1) 優先順位なし (有限バッファ、eager評価方式)

このスケジューラでは、データ、デマンドの到着に関わらず、マイクロプロセスを実行可能にし、レディ・キューに登録する。到着していないデータを待つ時、および、出力用のバッファが満たされている時に、待ち状態に移す。

(2) 優先順位付き (有限バッファ、eager評価方式)

これは、(1)に、優先順位を導入したものである。確実に必要な計算を行うマイクロプロセスに高い優先順位を与える。それ以外のマイクロプロセスも、実行可能になった段階で低い優先順位を与え、レディ・キューに登録する。

(3) 優先順位付き (有限バッファ、要求の先出し付き要求駆動)

このスケジューラでは、根に近い関数インスタンス (マイクロプロセス) から要求 (demand) を受け付けて初めて実行可能状態にする。並列性を抽出するために、要求の先だしを行う。これにより、確実に必要な計算を行うマイクロプロセスだけが実行可能となり、レディ・キューに登録される。

ここで、(3)により実現される方針が、SMASHシステムの本来のスケジューリング方針に近い。(2)では、SMASHにおける要求の発行時に優先順位を上げ、要求されたデータの生成の完了時に、優先順位を下げる。(1)、および、(2)では、結果的に無駄となる計算が行われることがある。(2)では、優先順位を利用することにより、無駄になる可能性が小さい計算を先に実行する。

ここで(1)を、あらかじめ用意してあるFIFOスケジューラを用いて実現した。これは、3.9.1項で述べたレベル5に対応する。また、(2)、(3)を、第2層 (図3-3) のレディ・キュー・モジュールを入替えることで実現した。これは、3.9.1項で述べたレベル4に対応する。

(2)では、マイクロプロセスの優先順位は、動的に決定される。これは、実時間システムで用いられる優先順位継承方式 (Priority Inheritance Protocols) の単純な場合になっている [74]。(c)の実現を、(b)において、優先順位が上がるまでマイクロプロセスの実行を遅延することにより行った。

ここで実現したスケジューラの性能を調べる実験を行った。その結果については、3.12.7項において延べる。



### 3. 1 0 仮想プロセッサを実現するカーネルの構成法

この節では、仮想プロセッサを実現するカーネルの構成法について述べる。この構成法の特徴は、利用者プロセスとカーネルの構造が相似である点にある。ここでは、まず仮想プロセッサの固有領域の実現について述べる。次に、これと相似の関係にある実プロセッサの固有領域について、その外部仕様、利用、および、内部実現について述べる。最後に、カーネル内部の軽量プロセス・スケジューラを用いたプロセスと仮想プロセッサのスケジューラの実現について述べる。

この節では、3. 6 節において述べた仮想プロセッサの外部仕様を例として、カーネルの構成法を述べる。しかしながら、ここで述べる方法は、汎用性があり、他の外部仕様を持つ仮想プロセッサやカーネル制御方式の軽量プロセスを実現する場合にも利用可能である。

#### 3. 1 0. 1 仮想プロセッサの固有領域の内部実現

この節では、3. 6 節で述べた仮想プロセッサの固有領域の外部仕様を実現する方法について述べる。論理的には、図3-10に示したようなアドレス変換

(図3-10の⑧と⑨)を実現するようにMMU (Memory Management Unit) の変換表 (translation table) を作成すればよい。しかしながら、単純に仮想プロセッサごと変換表を作成したならば、1つのプロセスについて仮想プロセッサの数だけの変換表が存在することになり、次の様な問題点が生じる。

(1) 変換表の操作の実行速度が低下する。たとえば、プロセスにメモリを割り当てる場合、そのプロセスに属するすべての仮想プロセッサについて変換表を書き換えなければならない。そのためメモリ資源を割り当てる操作の実行時間は、仮想プロセッサの数に比例することになる。

(2) 変換表のためのメモリ資源が無駄になる。

この項では、変換表の、利用者空間における共有領域に対応する部分を共有することにより上記の問題を解決する方法を示す。

##### 3. 1 0. 1. 1 2段階の変換表を利用するMMU

ここでは、図3-18に示すような2段階の変換表を利用するMMUにおける変換表の設定について述べる。最近のマイクロプロセッサの多くは、多段の変換表を持っており、ここで述べる方式を利用することができる [99]。たとえば、SPARCでは、3段階、MC68030では、4段階 (0段階から4段階まで可変)、インテル80386では2段階の変換表を利用する。実際に利用したMMUは、2段階の変換表を利用する、モトローラ社のMC88200である [52]。

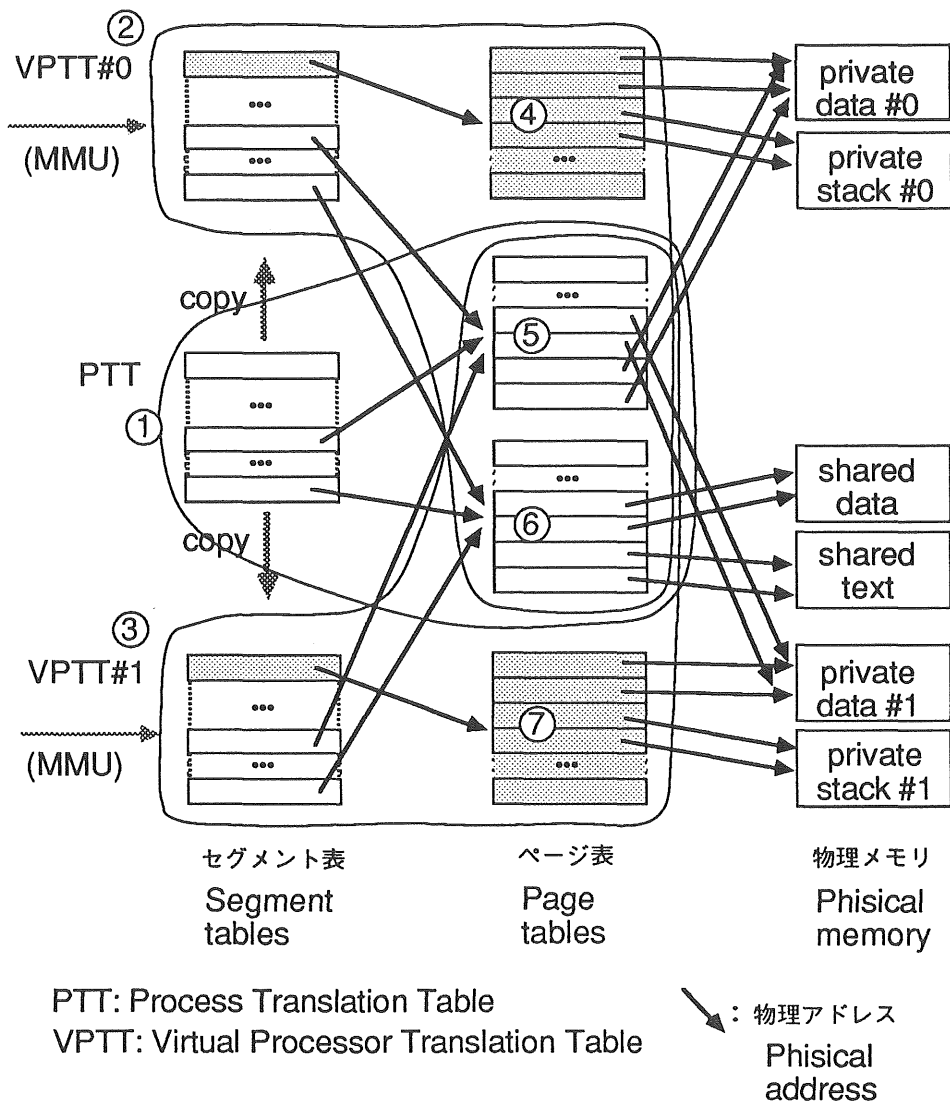


図 3 - 1 8 プロセス変換表と仮想プロセッサ変換表

Figure 3-18: Process and virtual processor translation tables

ここでは、2段階の変換表を用いることを仮定する。2段階以上の変換表の場合、ある段階で区切ることにより、2段階と同等に扱うことが可能である。ここでは、変換表の第1段をセグメント表 (segment table)、第2段をページ表 (page table) とよぶことにする。CPUが生成する論理アドレスは、上位ビットから

- ・セグメント番号
- ・ページ番号
- ・ページ内オフセット

に分けられる。アドレス変換は、次のようにして行われる。

- (1) 論理アドレス中のセグメント番号を添え字としてセグメント表のエントリを選択する。
- (2) そのエントリには、ページ表の物理アドレスが含まれている。それを用いて、ページ表を選択する。

(3) そのページ表において、論理アドレスのページ番号を添え字としてエントリを選択する。

(4) そのエントリには、メモリの物理アドレスが含まれている。このアドレスとページ内オフセットを加えて、物理アドレスを生成する。

### 3. 10. 1. 2 プロセス変換表と仮想プロセッサ変換表

この項では、プロセス用と仮想プロセッサ用の2種類の変換表を用いて仮想プロセッサのアドレス空間を設定する方法を示す。前者をプロセス変換表 (PTT, Process Translation Table, 図3-18の①), 後者を仮想プロセッサ変換表 (VPPT, Virtual Processor Translation Table, 図3-18の②と③) とよぶことにする。プロセス変換表は、共有領域の管理を行うためのものであり、MMUによって参照されることはない。仮想プロセッサ変換表は、その仮想プロセッサを実プロセッサが実行する時にMMUにより参照するための変換表であり、共有領域と固有領域の両方のアドレス変換に関する記述を含む。

図3-10に示した論理アドレスから物理メモリへの変換を実現する変換表を図3-18に示す。図3-18は、図3-10と同様に、アドレスが小さい(0に近い)方が下になるように描かれている。図3-10の⑧、⑨に示した変換を実現しているのは、図3-18の②、③に示す2つの仮想プロセッサ変換表である。それらの変換表(②、③)は、2つのページ表(⑤、⑥)を共有している。そして、物理ページも共有している。この部分は、3.6.3項において述べた共有領域に対応している。

図3-18において、影を付けた部分が固有領域に関係している部分である。たとえば、図3-10の⑤に示した固有データ・セグメントの論理アドレスは、②と④を用いて物理アドレスに変換される。

### 3. 10. 1. 3 共有領域へのメモリの割り付け

仮想プロセッサのアドレス空間の大部分は、共有領域である。したがって、共有領域への物理メモリの割り付けの速度が重要となる。ここでは、ページ表を共有している利点を活用し、高速化を図る手順を示す。

(1) プロセス変換表(図3-18の①)を根として、ページ表を共有しない方式と全く同じ手順を用いて、変換表を構築する。

(2) (1)において、プロセス変換表のセグメント表のエントリの中で変更したものを、そのプロセスに属するすべての仮想プロセッサの変換表のセグメント表(図3-18では②、③)にコピーする。(1)においてセグメント表を変更しなかった場合は、なにもしない。

(3) 実プロセッサが割り当てられている仮想プロセッサについて、MMU中の変換表のキャッシュ (TLB, Translation Look-aside Buffer、MC88200では、PATC (Page Table Translation Cache) と呼ばれている) をフラッシュする。

(4) 実プロセッサが割り当てられていない仮想プロセッサについて、変換表が変更されたことを示すマークを付ける。

これらの操作の中で、(2)の操作においてページ表を共有している利点が活かされている。もしページ表を共有していなかったならば、(1)の操作を仮想プロセッサの数だけ繰り返す必要が生じる。

ページ表を共有するためには、ある制限が必要である。それは、セグメント表の1エントリによって記述されるアドレス空間内に固有領域と共有領域の混在を許さないことである。この制限により、上の(2)において単純にセグメント表のエントリのコピー(浅いコピー、shallow copy)を行うだけでよいことになる。

### 3. 10. 2 実プロセッサの固有領域

3. 6. 3項と3. 6. 4項では、利用者プロセスにおける仮想プロセッサごとの固有領域の外部仕様と利用について述べた。前項(3. 10. 1項)では、仮想プロセッサごとの固有領域の内部実現について述べた。一方、カーネルの中においては、実プロセッサごとに必要となるデータが存在する。ゆえに、カーネルにおいて実プロセッサごとの固有領域を設定することは、有効である。この項では、その外部仕様、利用、および、内部実現について述べる。

#### 3. 10. 2. 1 実プロセッサの固有領域の外部仕様

実プロセッサの固有領域とは、実プロセッサのアドレス空間(カーネルの論理空間)の一部に存在する、実プロセッサごとに異なる物理メモリが参照される領域のことである。これに対して、実プロセッサの固有領域以外の領域を実プロセッサの共有領域とよぶことにする。

実プロセッサの固有領域は、利用者レベルのアドレス空間の設定から独立している。すなわち、他の外部仕様を持つ仮想プロセッサを提供するカーネルや、カーネル制御方式の軽量プロセスを提供するカーネルの実現においても、実プロセッサの固有領域を利用することができる。たとえば、Mac hシステムのカーネルを実現する場合においても、実プロセッサの固有領域を利用することが可能である。

#### 3. 10. 2. 2 実プロセッサの固有領域の利用

3. 6. 4項で行った議論は、仮想プロセッサを実プロセッサに置き換えてもそのまま成り立つ。

#### 3. 10. 2. 3 実プロセッサの固有領域の内部実現

実プロセッサの固有領域は、利用者空間とカーネル空間の変換表の形式が同じ場合、3. 10. 1項で述べた、仮想プロセッサの固有領域の実現とほとんど同じ技術を用いて実現することが可能である。これにより、利用者プログラムに機能を提供するためのモジュールとカーネル自身の動作のために必要なモジュールの共通化を図ること

ができるという利点が生じる。

利用者空間とカーネル空間の変換表の形式が異なる場合には、共通化を図ることができない。しかしながら、カーネル空間においてMMUが利用可能であるならば、実プロセッサの固有領域を実現することは、可能である。

固有領域の実現において、仮想プロセッサと実プロセッサで異なる操作が必要な点が存在する。これを以下にまとめる。

(1) カーネルには、固有スタック領域(スケジューラ用のスタック)が不要である。これは、自力で立上がる(bootstrap)時に利用するスタックをそのままスケジューラ用のスタックとして利用できるからである。

(2) カーネル用の変換表が完成するまで、セグメント表やページ表に用いるページの属性の設定を遅らせる必要がある。セグメント表やページ表に用いる物理メモリは、キャッシングの機能を停止しなければならない。そのためには、カーネル用のセグメント表やページ表があるページの属性を変更する必要がある。ところが、カーネル用の変換表を作成している途中の段階では、この属性変更を行うことは不可能である。したがって、1度変換表を作成した後に改めてセグメント表とページ表の物理アドレスを調べ、該当するページの属性を変更する必要がある。

### 3. 10. 3 仮想プロセッサとプロセスのスケジューラ

3. 10. 2項では、利用者レベルの抽象である仮想プロセッサの固有領域と同一の技術を用いて実プロセッサの固有領域を実現する方法について述べた。この項では、利用者レベルの軽量プロセスのスケジューラ(マイクロプロセス・スケジューラ)の技術を用いて、仮想プロセッサとプロセスのスケジューラを実現する方法について述べる。

#### 3. 10. 3. 1 プロセスと仮想プロセッサのスケジューラの外部仕様

プロセス・スケジューラ(大域スケジューラ(global scheduler))とは、資源割当ての単位であるプロセスの間で公平なCPU時間の分配を実現するためのスケジューラである。公平な分配を実現するために、実プロセッサの横取り(preemption)が行われる。

仮想プロセッサ・スケジューラとは、同一プロセス内において、別の仮想プロセッサへの制御の移動を行うためのスケジューラである。(本システムの仮想プロセッサ・スケジューラについては、3. 4. 2項において述べた)。制御の移動の原因としては、同期式入出力、同期式プロセス間通信、仮想記憶システムではページ・フォールト等に関する仮想プロセッサの停止/再開、および、制御の移動を要求するカーネル・コール(本システムでは、3. 6. 1項述べた `vp_switch()`)の発行があげられる。

ここでは、それぞれのスケジューラにおける方針(policy)とは独立した構成法を示す。ここで述べる構成法は、プロセスと仮想プロセッサという概念が存在するなら

ば、他の方針においても利用可能である。

### 3. 10. 3. 2 プロセス・スケジューラと仮想プロセッサ・スケジューラの内部構造

#### カーネル内の軽量プロセスとしての仮想プロセッサ

オペレーティング・システムのカーネルは、利用者からのトラップを扱う視点から、次の2つの方式に分類される [17]。

- (1) プロセス・モデル：利用者プロセスと1対1に対応し、固有のスタックを持つプロセスをカーネル内に設ける。後述する割込みモデルと比較して、実現が容易である。Mach 2.5 システム [2] [67]、UNIX システム [4] [46] において採用されている。
- (2) 割込みモデル：カーネルは、利用者からのトラップを割込みとして扱う。カーネル内には、プロセッサごとにスタックを設ける。前述のプロセス・モデルと比較して、効率的な実現が可能である。Mach 3.0 システム [17]、V システム [12] 等で採用されている。

本システムは、(1)を採用する。その理由は、仮想プロセッサをカーネル内の軽量プロセスとして実現することで、仮想プロセッサとプロセスのスケジューラをカーネル内の軽量プロセスのスケジューラとして実現可能であるからである。さらに、オペレーティング・システムのカーネルを1つの並列プログラムと考えた場合、(利用者レベルの抽象である)仮想プロセッサは、並列処理の単位となり得るため、軽量プロセスにより実現することが自然である。(カーネル中には、仮想プロセッサ以外にも軽量プロセスが存在する。)

#### カーネル内の軽量プロセスのグループとしての利用者プロセス

実プロセッサの制御という視点では、プロセスは、仮想プロセッサを実現している軽量プロセスのグループとしてとらえることができる。たとえば、プロセスを削除することは、そのグループに属するすべての(仮想プロセッサを実現している)軽量プロセスを削除することに相当する。プロセスを停止させることは、そのグループに属する実行中の(仮想プロセッサを実現している)軽量プロセスから、実プロセッサを奪いとることに相当する。

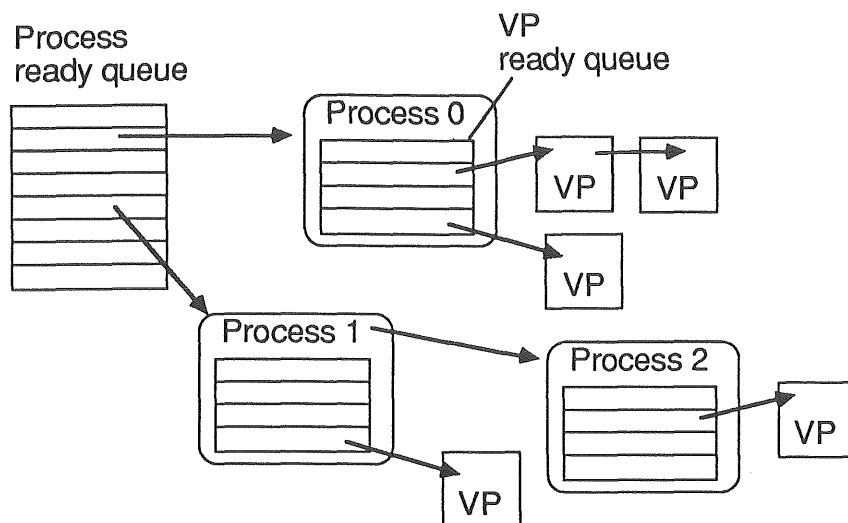


図 3 - 1 9 2 段階のレディ・キュー

Figure 3-19: Two-level ready queues.

### 2 段階のレディ・キューを持つスケジューラ

上で述べたように、仮想プロセッサは、カーネル内の軽量プロセスであり、プロセスは、カーネル内の軽量プロセスのグループである。この構造に従い、図 3 - 1 9 に示すような 2 段階のレディ・キューを用いてスケジューラを構築することができる。第 1 段は、(実行可能な仮想プロセッサを保持している) プロセスを保持するレディ・キューである。これを操作する部分が、3. 4. 2 節で述べたプロセス・スケジューラに相当する。第 2 段は、各プロセス内部の実行可能な仮想プロセッサを保持するレディ・キューである。これを操作する部分が、3. 4. 2 節で述べた仮想プロセッサ・スケジューラに相当する。

図 3 - 1 9 は、プロセスのレベルにおいて 8 つ、仮想プロセッサのレベルにおいて 4 つの優先順位に分割されたレディ・キューを示している。実プロセッサは、プロセス・レディ・キューからプロセスを選ぶ。そしてそのプロセス中の仮想プロセッサ・レディ・キューから実行可能な仮想プロセッサを得て、それを実行する。

### プロセスの操作

このような構造を持つスケジューラでは、プロセスの操作は、以下のように実現される。

(1) 停止：停止するプロセスをプロセス・レディ・キューから削除する。そのプロセスの実行中状態の仮想プロセッサについて、それを実行している実プロセッサに割り込みを行う。割り込まれた実プロセッサは、実行していた仮想プロセッサを、そのプロセスの仮想プロセッサ・レディ・キューに入れる。実プロセッサは、他のプロセスの仮想プロセッサに制御を移す。

(2) 実行再開：実行を再開するプロセスを単にプロセス・レディ・キューに接続

するだけでよい。

(3) 生成：アドレス空間や最初の仮想プロセッサを生成した後、(2)を行う。

(4) 終了、および、削除：(1)の操作を行った後、プロセスが保持している資源を開放し、終了コードを保存する。

このように、プロセスの操作は、プロセス・レディ・キューへの操作により容易に実現される。

### 3. 1 1 実現

3. 3節から3. 1 0節において述べたマイクロプロセス／仮想プロセッサ方式に基づき、軽量プロセス機能を実現した。この節では、はじめに、裸の共有メモリ型マルチプロセッサにおける仮想プロセッサを提供するカーネルの実現について述べる。次に、カーネル外のマイクロプロセス・ライブラリの実現について述べる。最後に、既存のシステムにおける仮想プロセッサ機能のエミュレーションについて述べる。

#### 3. 1 1. 1 仮想プロセッサを実現するカーネルの実現

##### 3. 1 1. 1. 1 実現環境

3. 1 0節で述べた構成法に基づき、共有メモリ型マルチプロセッサ Luna 8 8 k [ 3 2 ] [ 1 0 5 ] において仮想プロセッサを提供するカーネルの実現を行った。実現に用いたLuna 8 8 kは、4プロセッサ構成で、主記憶容量は、3 2 Mバイトである。各プロセッサは、命令とデータのそれぞれに1 6 kバイトのキャッシュを供えている。Luna 8 8 kに付属のオペレーティング・システムであるUnioS - Mach 1. 1 0 (Mach 2. 5相当)を開発システムとして利用した。

##### 3. 1 1. 1. 2 モジュール構成

カーネルの開発には、C言語を用いた。全体の大きさは、約2万1 0 0 0行である。モジュール構成を、以下に示す。

#### メモリ管理

bitmap	メモリ管理用ビットマップ操作モジュール	5 8 5 行
mo	メモリ・オブジェクト管理	4 1 1 行

#### ハードウェア操作

cmmu	CMMU M88200操作モジュール	2 5 1 5 行
dev	デバイス管理	1 7 8 0 行



fp	浮動小数点用カーネル内ルーチン	2 5 9 0 行
m88k	MPU M88100関連手続き	1 4 7 7 行
env	コンテキスト切り替え	2 1 9 行

#### 割込み受付け

locore	割込み受付けモジュール	2 7 1 5 行
main	カーネル・コール受付け	2 6 4 0 行

#### プロセスと仮想プロセッサの実現

mp	カーネル内マイクロプロセス	7 8 6 行
vp	仮想プロセッサ操作	1 1 1 0 行
ps	プロセス・オブジェクト操作	1 2 6 6 行
rp	実プロセッサ操作 (スケジューラ)	8 7 8 行

#### カーネル・サーバ呼び出し受付け部

con	コンソール・ファイル・サーバ	1 5 8 行
time	時間サーバ	8 4 行
mem	メモリ・サーバ	5 1 行
stdps	プロセス・サーバ	5 4 0 行

#### その他

misc	その他ユーティリティ手続き	1 5 8 8 行
------	---------------	-----------

---

合計	2 0 9 0 7 行
----	-------------

これらのモジュールの中で、浮動小数点用カーネル内ルーチン (fp) と割込み受付けモジュール (locore) については、オムロン社より提供されたプログラムを変更して利用した。デバイス管理については、オムロン社より提供されたプログラムを参考にして開発を行った。

### 3. 1 1. 1. 3 カーネル利用者用ライブラリ

カーネル利用者用ライブラリとは、上で述べた L u n a 8 8 k 上の R e S C カーネルを利用してプログラムを開発するために C 言語から利用可能な関数を集めたものである。これには、上記のカーネルを呼び出すためにトラップ命令を含むものと、一般の C 言語用のライブラリが含まれる。

services	インタフェース記述	2 1 6 行
----------	-----------	---------

合計

912行

上記の他に、カーネルと共通のルーチンがある。さらに、開発システム上のC言語のライブラリ関数から、UNIXのシステム・コールを含まないものを抜出して利用している。このような関数の例としては、文字列操作ライブラリ関数群 (string(3)) があげられる。

### 3.11.1.4 考察

3.10.1項、および、3.10.2項において述べた仮想プロセッサと実プロセッサの固有領域の実現においては、ハードウェア (MMU) を操作する部分を除き、その大部分を開発システム上で利用者プログラムとしてデバッグと動作の確認を行うことができた。この時、ハードウェアを操作する手続きについては、ハードウェアのレジスタに相当するダミーの変数を設け、その内容を確認することにより、動作の確認を行った。また、一部の手続きについては、ハードウェアのエミュレートを行うことにより、デバッグを行った。エミュレートを行った手続きの例きとしては、CMMUのアドレス変換操作や文字入出力があげられる。

3.10.3項において述べたスケジューラについても、割込み制御の部分を除き、その大部分を開発システム上でデバッグと動作の確認を行うことができた。これらの実現を通じて、カーネルを利用者プログラムと同様に開発することが可能であることを確認することができた。

今回の実現には、ディスク入出力を行うファイル・サーバの機能が含まれていない。このファイル・サーバの機能は、2.10.4項で述べたように、NFSサーバ [59] 中のファイル・サーバの機能からキャッシングの機能を取り外したようなものになる。(NFSサーバは、他に名前サーバとしての機能を持つ。) したがって、NFSサーバの一部を抜き出すことにより、ファイル・サーバを開発することが可能であると思われる。

## 3.11.2 マイクロプロセス・ライブラリの実現

### 3.11.2.1 実現環境

3.5章で述べたマイクロプロセス・ライブラリを、5種類のソフトウェア環境 (前項で述べた本RES Cシステムのカーネルを含む) と4種類のハードウェア・アーキテクチャにおいて実現を行った。すなわち、5種類の第0層と4種類の第1層を開発した。第2層以上は、これらのソフトウェア環境とハードウェア・アーキテクチャから独立している。

第0層の実現を行ったソフトウェア環境を以下に示す。

- (1) ReSCカーネル (本システム)
- (2) Dynix (Sequent Balance 8000)
- (3) SunOS (Sun3, Sun4)
- (4) Mach2.5 (Luna88k, NeXT)
- (5) NewsOS (News)

(1) については、3.11.1で述べたように、トラップ命令を含む手続きをアセンブリ言語により記述した。その他については、UNIXの機能を利用してReSCカーネルの機能のエミュレーションを行った。エミュレーションの詳細については、3.11.3項において述べる。

第1層の実現を行ったハードウェア・アーキテクチャを以下に示す。

- (1) MC88100 (Luna88k)
- (2) NS32032 (Sequent Balance 8000)
- (3) MC68000 (Sun3, NeXT, News)
- (4) SPARC (Sun4)

### 3.11.2.2 モジュール構成

マイクロプロセス・ライブラリは、3.5章において述べたように、4つの層から構成されている。それぞれの層におけるモジュールの名前、機能、および、大きさは、以下のようになっている。

#### 第0層

resc	カーネル利用者ライブラリ	(3.11.1項)
	ReSCカーネルのエミュレーション	1100行
mem	メモリ管理	3764行

#### 第1層

env	コンテキスト切り替えルーチン	868行
lock	スピンロック・ルーチン	418行

#### 第2層

main	生成消滅と同期、初期化	1660行
trace	実行の軌跡の記録と解析	1025行
misc	その他ユーティリティ手続き	693行

### 第3層

mon	モニタ (1)	2 2 4 行
mutex	モニタ (2)	2 1 8 行
sem	セマフォ	1 2 2 行
msg	ランデブ	3 0 6 行
pipe	ストリーム	8 7 7 行

### 合計

---

合計	1 1 2 7 5 行
----	-------------

### ■第0層

カーネル利用者ライブラリについては、3. 1 1. 1. 3 において述べた。R e S Cカーネルのエミュレーションに関しては、1 1. 3 節において述べる。メモリ管理 (mem) は、仮想プロセッサの共有領域の実現と汎用メモリ管理ルーチン (U N I X における C 言語のライブラリ関数 malloc(3)、free(3) に対応) を含んでいる。なお、後者については、GNU E m a c s に含まれていたものにスピロックによる相互排除機能を付加することにより、実現を行った。

### ■第1層

上で示したモジュールの行数は、3. 1 1. 2. 1 において述べた4種類のアーキテクチャの合計である。コンテキスト切り替えについては、アセンブリ言語により記述を行った。主に、レジスタの退避、スタックの切替え、レジスタの回復から構成されている。スピロックについては、原子的な命令を用いて実現を行っている。ただし、図3-9に示したものとは異なり、ロックの解除を待つ部分を、通常のロード命令を用いて実現している。このため、キャッシュの効果により、共有バスへのトラヒックが抑えられる。

### ■第2層

mainが、第2層の主要部である。trace は、3. 5. 3 項で述べたマイクロプロセスの実行の軌跡の記録と解析を行うモジュールである。最後のmiscには、第2層を動作させるためのユーティリティが含まれている。たとえば、キューの操作、メッセージの表示、シンボル・テーブルの解析、環境変数の解析を行う関数が含まれている。

### ■第3層

ここで、モニタ (1) は、高レベルの、モニタ (2) は、低レベルの機能を使ったモニタ同期プリミティブの実現である。前者は、第2層において提供するデッドロックを検出する機能が有効であるのに対して、後者は無効である。また、前者は、

SunLWPにおけるモニタと、後者は、C-Threadsにおけるmutexと似た内部構造を持つと思われる。

### 3.11.2.3 考察

この実現においては、層化プログラミング (layered programming) の技術により、移植性と可読性が高いライブラリを構築することができた。カーネルに依存する部分は、第0層に、そしてプロセッサのアーキテクチャに依存する部分は、第1層に隠蔽されている。第3層の各同期通信プリミティブの構築が非常に簡単になっていることから、第2層において提供した機能の妥当性を確認することができた。

### 3.11.3 UNIXにおける仮想プロセッサのエミュレーション

本ReSCシステム以外においてもマイクロプロセス・ライブラリを利用するために、UNIXのプロセスを用いて仮想プロセッサのエミュレーションを行った。実現を行ったのは、Dynixシステム [6]、Mach 2.5システム [105] [58]、SunOS [93] である。

基本的には、UNIXのプロセスを仮想プロセッサとして利用し、そしてマイクロプロセス用のメモリを、UNIXのプロセス間共有メモリにより実現する方法を用いた。これは、文献 [106] で用いられている方法と同じ方法である。

#### 3.11.3.1 カーネル・コール

vp\_allocate() カーネル・コールを、UNIXの fork() システム・コールにより実現した。vp\_sleep() と vp\_switch() カーネル・コールを、UNIXの usleep() ライブラリ関数により実現した。(usleep()は、ソフトウェア割込みにより実現されている。)

#### 3.11.2.2 仮想プロセッサの固有領域と共有領域

Dynixでは、C言語により提供されている private/shared キーワード、および、shm\_alloc() ライブラリ関数を用いて仮想プロセッサの固有領域と共有領域を実現した。内部的には、それらは、UNIXの mmap() システム・コールをインタフェースとして実現されている。

Machシステム、および、SunOSでは、固有領域を通常のUNIXのデータ・セグメントにより実現した。共有領域については、それぞれ vm\_allocate() カーネル、および、mmap() システム・コールを用いて、プロセス間共有メモリを実現し、その上に共有領域を割り当てた。共有領域上の静的変数については、データ・セグメントからある定数のオフセットを加えた領域にプロセス間共有メモリを割り当て、その部分を参照することで実現した。

### 3. 1 1. 2. 3 資源割当て

プロセスを用いて仮想プロセッサを実現した場合、もっとも大きな問題は、資源の割当てである。たとえば、ある1つの仮想プロセッサにおいて、ファイルを開いたとしても、他の仮想プロセッサには伝わらない。ファイルを閉じる操作についても同様である。この問題は、エミュレーションに起因するものである。本R e S Cシステムでは、カーネルにおいて資源割当ての単位としてのプロセスと実プロセッサのエントリとしての仮想プロセッサを実現しているため、このような問題は生じない。

この問題を解決するために、U N I Xのシステム・コール `sendmsg()`、`recvmsg()`を用いて、ファイル記述子（資源へのポインタ）を他のプロセスへ転送する方法を用いた。以下に、ファイルを開く操作（`open`システム・コール）を例に、手順を説明する。

#### 準備：

仮想プロセッサごとに、双方向の通信が可能なソケットを生成しておく。

#### ファイルを開く仮想プロセッサ：

- (1) U N I Xの `open()` システム・コールを発行し、ファイルを開く。エラーが生じた時には、これを返す。
- (2) 自分以外の仮想プロセッサ（を実現しているプロセス）のソケットへ `sendmsg()` システム・コールを利用して、ファイル記述子を送る。
- (3) 自分以外の仮想プロセッサ（を実現しているプロセス）へ、`kill()` システム・コールを用いて、ソフトウェア割込みを行う。
- (4) `recvmsg()` システム・コールを用いて、自分以外の仮想プロセッサからの応答を待つ。

#### それ以外の仮想プロセッサ：

- (1) それ以外の仮想プロセッサ（を実現しているプロセス）は、ソフトウェア割込みを受け付けると、自分自身のソケットから `recvmsg()` システム・コールを用いて、メッセージを受け取る。これにより、ファイル記述子が渡されることになる。
- (2) ソフトウェア割込みを行った仮想プロセッサに対して、`sendmsg()` システム・コールにより、応答を送る。

以上は、ファイルを開く操作の手順であった。その他に、システム・コール `creat()`、`pipe()`、`dup()`、`dup2()`、`socket()` が、ほぼ同様の手順により実現される。システム・コール `close()` の場合、ファイル記述子ではなく、閉じるべき番号を整数として送り、ソフトウェア割込みを受付けた仮想プロセッサにおいて、`close()` システム・コールを発行する。

### 3. 1 1. 2. 4 限界

UNIXのシステム・コールを利用したReSCカーネルのエミュレーションには、限界がある。第1に、最も深刻な問題は、プロセス全体の操作である。たとえば、プロセスを終了させる場合、仮想プロセッサを構成している全てのUNIXのプロセスを終了させる必要がある。この時、終了させることに失敗したり、あるいは、個々のUNIXのプロセスが勝手に終了する危険性がある。エミュレーションを行なっている限り、この問題を解決することは、不可能である。

第2に、メモリ資源割当ての問題がある。1度並列処理を開始した後(fork()システム・コールを発行して、複数の仮想プロセッサを生成した後)では、利用するメモリ資源の量を拡張することができない。この問題は、仮想記憶を利用することで、ある程度緩和される。

第3に、Machシステム、および、SunOSでは、仮想プロセッサ間の共有領域に対応するコンパイラ、アセンブラ、および、リンケージ・エディタが存在しない。したがって、これらのシステムでは、メモリの無駄が生じている。

最後に、マイクロプロセス・スケジューラが働かなくなるという問題がある。マイクロプロセス・スケジューラは、特定のマイクロプロセスに多くの資源を与えることで、その並列応用プログラム全体の実行時間の短縮を試みる。一方、カーネルは、仮想プロセッサを実現しているプロセスを公平に実行しようとする。このことは、1つの並列応用プログラムの実行において、利用者レベルとカーネル・レベルの2つのスケジューラが動作していることを意味する。その結果、利用者レベルのマイクロプロセス・スケジューラが働かなくなる。

### 3. 1 2 性能

この節では、実現した軽量プロセス機能の性能について述べる。

#### 3. 1 2. 1 実験環境

実現したマイクロプロセス、および、仮想プロセッサの性能を調べる実験を行った。実験に用いたシステムのプロセッサのアーキテクチャ、既存のオペレーティング・システム、主記憶容量、キャッシュ容量（制御方式）を以下に示す。

#### 共有メモリ型マルチプロセッサ

- (1) S e q u e n t   B a l a n c e 8 0 0 0、  
D y n i x 2 . 1、  
4 × N S 3 2 0 3 2   1 0 M H z、  
8 M   B y t e s、  
8 k   B y t e s ( w r i t e - t h r o u g h )
  
- (2) O m r o n   L u n a 8 8 k、  
M a c h 2 . 5   ( U n i O S - M a c h   1 . 1 0 )、  
4 × M 8 8 1 0 0   2 5 M H z、  
3 2 M   B y t e s、  
3 2 k   B y t e s   ( c o p y - b a c k )

#### 単一プロセッサ

- (3) S P A R C s t a t i o n   2、  
S u n O S   4 . 1 . 1、  
S P A R C   4 0 M H z、  
3 2 M   B y t e s、  
6 4 k   B y t e s
  
- (4) S u n 3 / 6 0 (以下 S u n 3 と略記する)、  
S u n O S   4 . 0 . 3、  
M 6 8 0 2 0   1 6 M H z、  
1 2 M   B y t e s



(5) NeXT (cube)、  
Mach 2.5、  
M68030 25MHz、  
8M Bytes、  
on-chip

このように、様々な環境において実験を行うことにより、軽量プロセスの実現方式の間の、プロセッサのアーキテクチャや具体的な実現から独立した比較が可能になる。

(1)と(2)は、共に4プロセッサ構成の共有メモリ型マルチプロセッサであるが、プロセッサのアーキテクチャが大きく異なる。(1)は、CISC (Complex Instruction Set Computer) 系のプロセッサを、(2)は、RISC (Reduced Instruction Set Computer) 系のプロセッサを備えている。RISCの場合、CISCと比較してキャッシュが性能に大きく影響すると予想される。

(3)と(4)では、同一のオペレーティング・システムが動作しているが、プロセッサのアーキテクチャが大きく異なる。(3)は、RISC系のプロセッサを、(4)は、CISC系のプロセッサを備えている。特に(3)のSPARCプロセッサは、巨大なレジスタ・ファイルを持っていることから、軽量プロセスのコンテキスト切替えの処理時間が長くなることが予想される。SPARCプロセッサでは、手続き1個あたり16個のレジスタを局所変数として利用することができる。コンテキスト切替えにおいては、本来、利用しているレジスタのみを保存/回復すれば十分である。しかしながら、レジスタの利用状況を得る手段がないため、全てのレジスタを保存/回復しなければならない。さらに、レジスタ・ファイルのフラッシュを特権モードで行う必要があるため、コンテキスト切替えにおいて必ずカーネル・コールが必要となる。

(5)と(2)では、同一のオペレーティング・システムが動作しているが、プロセッサのアーキテクチャと数が異なる。

以下の3.12.3項、3.12.4項、3.12.5項、および、3.12.7項の実験では、既存のシステム上に実現したマイクロプロセス・ライブラリを用いて実行時間を測定した。3.12.6項の実験では、裸の計算機上に実現したカーネルの性能を測定した。

以下の実験においては、実行時間を、UNIXのシステム・コール `gettimeofday()`、あるいは、ReSCカーネルの同等の機能を用いて測定した。実行時間の精度が異なるのは、システムごとに測定可能な最小時間が異なること、および、繰り返し回数が異なることによる。既存のシステムでは、マルチユーザ・モード(いくつかのサーバ・プロセス (UNIXでは、デーモン (daemon)) が生成されている状態)において、システムの負荷が低い状態で実験を行った。Luna 88k上のReSCカーネルでは、測定対象となるプロセスのみを生成して実験を行った。

既存のシステム上での実験では、あらかじめメモリを参照することで、測定時間に

仮想記憶の処理時間が含まれないようにした。(Luna 88k上のResCカーネルは、実記憶システムであり、仮想記憶の影響はない。)さらに、プロセス生成・消滅の実験では、仮想記憶の影響により実行時間が長くない範囲内において、生成する軽量プロセスの数を大きくした。以下で示す実験結果では、1回あたりの実行時間に加えて、反復の回数を示す。これらの結果は、並列処理に用いることが可能な軽量プロセス数や同期処理の頻度の目安を与える。本方式では、軽量プロセスの数の上限は、利用者プロセスのアドレス空間の大きさによってのみ制約される。カーネル制御方式では、しばしばカーネル内に静的に確保された表の大きさにより、制約を受ける。

### 3.12.2 比較対象

以下の3.12.3項、3.12.4項、および、3.12.5項では、マイクロプロセスの基本的な性能を論じる。比較対象としては、カーネル制御方式のMachのC-Threads [105] [58]、コルーチン方式のSunOSのSunLWP [93]、および、重量プロセスのUNIXのプロセス [46]を取り上げる。

C-Threadsは、Luna 88k、および、NeXTシステムで利用可能である。SunLWPは、SPARCstation 2、および、Sun3で利用可能である。Sun3の実行形式の一部は、NeXTのatomコマンドを利用することにより、NeXTシステム用の実行形式へ変換することが可能である。この機能を利用して、SunLWPをNeXTシステム上でも動作させ、実行時間を測定した。なお、Balance 8000では、他の軽量プロセスを利用することができなかったため、UNIXの重量プロセスとの比較だけを行った。

#### ■ C-Threads

C-Threadsでは、軽量プロセスの生成・消滅、および、コンテキスト切替えには、カーネル・コールが必要である。コンテキスト切替えを伴わない同期は、カーネル・コールを用いることなく実現される。

C-Threadsでは、軽量プロセスのキャッシングを行っている。すなわち、終了した軽量プロセスを破壊せずにそのまま保存しておき、次に生成要求を受け付けた時にそれを再利用している。

この実験において用いたC-Threadsは、Machシステムの第2.5版に付属のものであり、C言語レベルのCスレッドとカーネル・レベルのスレッドの間には、1対1の対応関係が存在する。第3版に付属のものでは、両者の間に1対1の対応関係が存在しない。

#### ■ SunLWP

SunLWPは、SunOS 4.xに付属している軽量プロセス・ライブラリであ

る。これは、コルーチン方式により実現されている。軽量プロセスのスタックをレッド・ゾーン方式により保護する機能がある。この機能は、SunOS固有のシステム・コール `mprotect()` により実現されている。SunLWPでは、そのスタックのキャッシングを行う機能がある。

SPARCstation 2では、コンテキスト切換えにおいて、レジスタ・ウィンドウをフラッシュするために、カーネル・コールを発行する。本方式においても、同一のカーネル・コールを利用した。

## ■UNIXのプロセス

従来の重量プロセスとの性能を比較するために、UNIXのプロセスを利用した。プロセスの生成のプリミティブとしては、UNIXの `fork()` システム・コールを利用した。プロセス間の同期のプリミティブとしては、セマフォ関係のシステム・コールを利用した。なお、Luna88kとNeXTでは、セマフォが利用できなかったので、比較を行わなかった。

### 3. 1 2. 3 プロセスの生成・消滅

#### ■性能測定用プログラム

2重ループ構造をもつプログラムを用いて、プロセスの生成・消滅操作の性能を測定した。内側のループでは、プロセスの入れ子を作り、共有メモリ型マルチプロセッサにおける並列処理の効果、および、軽量プロセスやスタックのキャッシングの効果を抑えている。逆に、外側のループの実行回数を変化させることで、キャッシングの効果を調べることができる。このプログラムを用いて、プロセスの生成処理、コンテキスト切替え（親から子へ、子から親へ）、同期処理（子の終了待ち）、消滅処理の全体の処理時間を測定した。付録D-1に、使用した性能測定用のプログラムの骨格を示す。

表3-1(a)と表3-1(b)に、実験の結果を示す。これらの表で、上段が実行時間、下段が生成したプロセスの数を表す。下段の数字が

$$m \times n$$

とは、外側のループを  $m$  回、内側のループを  $n$  回実行したことを意味する。

表3-1(b)における実行時間を、次の式により計算した。

$$(m + 1) \times n \text{ 回の実行時間} - 1 \times n \text{ 回の実行時間}$$

---

$$m \times n$$

表 3 - 1 プロセスの生成・消滅の実行時間の比較

上段：プロセス 1 個当りの実行時間（単位：ミリ秒）

下段：繰返し回数（単位：回）

Table 3-1: Comparisons of execution times of process creation and termination.

Upper: Execution times of process creation and termination. (In milliseconds)

Lower: The numbers of iteratins.

表 3 - 1 ( a ) 1 回目

Table 3-1(a): In the first creation of processes.

	SPARCstation 2	Sun3	NeXT	Luna88k	Balance 8000
Microprocess (Layer 2)	0.149 (1x1k)	0.400 (1x500)	0.444 (1x500)	0.19 (1x1k)	1.5 (1x100)
SunLWP (mprot)	0.97 (1x100)	3.2 (1x100)	— —	— —	— —
SunLWP (none)	0.43 (1x100)	1.0 (1x100)	1.03 (1x100)	— —	— —
C-Thread	— —	— —	10.0 (1x10)	10.5 (1x20)	— —
UNIX (heavy-weight)	6.60 (1x50)	16. (1x20)	27.7 (1x20)	16. (1x20)	38. (1x10)

■他の軽量プロセスとの比較

表 3 - 1 ( a ) からわかるように、スタックの保護を行った SunLWP(mprot)と比較すると、本実現 (Microprocess) が 10 倍から 20 倍高速になっている。主な原因は、SunLWPでは、レッド・ゾーン方式によりスタックの保護を行っていることにある。SunLWPにおいて、スタックの保護機能を外したとしても (表 3 - 1 ( b ) の SunLWP(none))、本実現は、約 3 倍高速になっている。この差は、軽量プロセスの実現方式の違いによるものではなく、両実現の機能の違いによる。たとえば、SunLWPでは、UNIX固有のソフトウェア割込みや例外を扱う機能があり、本マイクロプロセスよりも複雑になっている。軽量プロセスの実現方式の違いだけならば、3. 4. 1 項で述べたように、コルーチン方式の SunLWP の方の効率がよくなる。

表 3 - 1 プロセスの生成・消滅の実行時間の比較

上段：プロセス 1 個当りの実行時間（単位：ミリ秒）

下段：繰返し回数（単位：回）

Table 3-1: Comparisons of execution times of process creation and termination.

Upper: Execution times of process creation and termination. (In milliseconds)

Lower: The numbers of iterations.

表 3 - 1 ( b ) 2 回目以降

Table 3-1(b): In the second creation of processes.

	SPARCstation 2	Sun3	NeXT	Luna88k	Balance 8000
Microprocess (Layer 2)	0.149 6 ([11-1]x1k)	0.384 ([11-1]x500)	0.362 ([11-1]x500)	0.130 ([11-1]x1k)	1.49 ([11-1]x100)
SunLWP (mprot)	0.443 ([11-1]x100)	1.46 ([11-1]x100)	— —	— —	— —
SunLWP (none)	0.339 ([11-1]x100)	0.98 ([11-1]x100)	1.01 ([11-1]x100)	— —	— —
C-Thread	— —	— —	2.37 ([11-1]x10)	1.9 ([11-1]x20)	— —
UNIX (heavy-weight)	6.56 ([11-1]x50)	16.2 ([11-1]x20)	28.11 ([11-1]x20)	16.25 ([11-1]x20)	38. ([11-1]x10)

C-Threads と比較すると、キャッシングが有効でない場合（表 3 - 1 ( a ) ）、本実現が約 2 0 から 5 0 倍高速になっている。この差は、軽量プロセスの実現方式の違いによる。本方式では、カーネル・コールのオーバーヘッドがなく、必要となるデータ構造を全て利用者空間に置くことができるので、高速になっている。

また、表 3 - 1 ( a ) と表 3 - 1 ( b ) を比較すると、C-Threads、および、SunLWP (mprot) において、2 回目以降の実行時間が 1 回目と比較して著しく短縮していることがわかる。これは、C-Threads では、軽量プロセスのキャッシング、SunLWP (mprot) では、スタックのキャッシングが有効に働いているからである。本実現では、軽量プロセスやスタックのキャッシングを行っていないが、メモリのキャッシュの効果により、処理時間が短くなっている。

## ■UNIXのプロセスとの比較

表3-1において、マイクロプロセスとUNIXの重量プロセス（UNIX heavy-weight）を比較すると、マイクロプロセスが25倍から80倍高速であることがわかる。表3-1（a）と（b）を比較すると、UNIXでは、プロセスのキャッシングが行われていないことがわかる。

### 3.12.4 コンテキスト切替えを伴わないプロセス間の同期

#### ■性能測定用プログラム

コンテキスト切り替えを伴わないプロセス間の同期操作の性能を、次の様な同期プリミティブを繰り返し実行することにより測定した。

（1）（スピン）ロックのロック操作／アンロック操作

本マイクロプロセス・ライブラリの第1層（図3-3）で利用可能である。

（2）セマフォのP命令／V命令

本マイクロプロセス・ライブラリの第3層（図3-3）、  
および、UNIXで利用可能である。

（3）モニタの入口操作／出口操作

本マイクロプロセス・ライブラリの第3層（図3-3）、SunLWP、  
および、C-Threadsで利用可能である。

実験では、1つのプロセスにおいて、これらの操作を繰り返し行い、実行時間を測定した。実験の結果を、表3-2に示す。これらの表の上段は、同期操作1回当たりの実行時間、下段は、その実行時間を測定した時に行った操作の回数である。ここでは、ロックとアンロック、P命令とV命令、入口操作、出口操作を合せて1回の操作と数える。

表3-2において、mon(1)とmon(2)は、共にモニタの実現であるが、機能と実現が異なる。前者では、第2層で提供されるデッドロックを検出する機能が有効であるのに対して、後者では無効である。それは、後者では、スピンロックにより、モニタ入口手続きを実現しており、スピンロックが保持できなかった場合、待ち状態にするのではなく、mp\_yield()により自分自身を実行可能状態に置いたまま、他のマイクロプロセスへ制御を移しているからである。このような実現は、単一プロセッサでは、有効である。また、前者は、内部的な実現がSunLWPのものに近いと思われる。後者は、C-Threadsのものに近くなっている。

なお、第2層には、コンテキスト切替えを伴わないプロセス間の同期を実現するプリミティブが存在しないので、実行時間を測定していない。この層では、必要ならば、第1層のスピンロックがそのまま利用される。

測定に用いたプログラムの骨格を、付録D-2に示す。

表 3 - 2 コンテキスト切替えを伴わないプロセス間の同期の実行時間の比較

上段：同期操作の実行時間（単位：μ秒）。

下段：繰返し回数（単位：回）。

Table 3-2: Comparisons of execution times of inter-process synchronization without context switches.

Upper: Execution times of synchronization operations. (In microseconds)

Lower: The numbers of iterations.

	SPARCstation 2	Sun3	NeXT	Luna88k	Balance 8000
Microprocess Layer 3/mon(1)	27.1 (100k)	31.2 (100k)	29.9 (10k)	8.3 (100k)	232. (10k)
Microprocess Layer 3/mon(2)	0.672 (1000k)	6.2 (100k)	5.48 (100k)	2.22 (1000k)	57. (10k)
Microprocess Layer 3/sem	3.13 (100k)	34. (10k)	27.1 (10k)	9.7 (1000k)	150. (10k)
Microprocess Layer 1/spin	0.672 (1000k)	4.8 (100k)	3.92 (100k)	2.05 (1000k)	54. (10k)
SunLWP (monitor)	5.6 (100k)	40.8 (100k)	39.9 (10k)	— —	— —
C-Thread (mutex)	— —	— —	4.3 (100k)	2.1 (100k)	— —
UNIX (semaphore)	221. (10k)	1120 (1k)	— —	— —	1930. (1k)
UNIX (getpid())x2)	31.84 (100k)	196. (10k)	158. (1k)	68.4 (100k)	530. (10k)

mon: monitor  
sem: semaphore  
spin: spin lock

### ■ マイクロプロセスの基本的な性能

本マイクロプロセス・ライブラリの各層の機能を比較すると、第1層のスピンロック Layer 1/spin と第3層のモニタ mon(2) の性能がよいことがわかる。セマフォとモニタ mon(1) は、SPARCstation 2 を除いて同程度のオーバーヘッドになっている。SPARCstation 2 において、このように大きな差が生じるのは、レジスタ・ファイル

のフラッシュに関係している。すなわち、本ライブラリが層構造を持っているため、層に対応した手続き呼出しのネストが深くなっているからである。

#### ■他の軽量プロセスとの比較

表3-2において、本実現 Microprocess Layer 3/mon(1) と SunLWP(monitor) と比較すると、SPARCstation 2 では、本実現が5倍程度遅くなっているのに対して、逆に、Sun3 と NeXTでは、本方式が25%程度高速になっている。この差は、軽量プロセスの実現方式の違いによるものではなく、両者の機能と内部実現の違いによる。実現方式としては、相互排除のオーバーヘッドが不用なコルーチン方式の方が効率がよくなる。

表3-2において、本実現 Microprocess Layer 3/mon(2) と C-Threads(mutex)を比較すると、ほぼ同程度の性能であることがわかる。これは、3.5.2項で述べたように、カーネル制御方式の C-Threads においても、コンテキスト切替えを伴わないプロセス間の同期においては、カーネル・コールを発行することなく実現することが可能になっているからである。

#### ■UNIXのプロセスとの比較

表3-2において、本実現 Microprocess Layer 3/sem と、UNIX semaphore を比較すると、本実現が10倍から70倍程度高速になっている。UNIXにおける最も軽いシステム・コール getpid() と比較しても、さらに高速になっている。このことは、カーネル・コールを減らす事が高速処理につながることを示している。

### 3.12.5 コンテキスト切替えを伴うプロセス間の同期

#### ■性能測定用プログラム

コンテキスト切り替えを伴うプロセス間の同期操作の性能を、次の様な同期プリミティブを繰り返し実行することにより測定する。

- (1) 2つの私有セマフォを用いる。

本マイクロプロセス・ライブラリの第3層、UNIXにおいて利用可能である。

- (2) 1つのモニタと2つの条件変数を用いる。

本マイクロプロセス・ライブラリの第3層、C-Threads、および、SunLWPにおいて利用可能である。

- (3) コンテキスト切り替えを行う機能と呼出す。

本マイクロプロセス・ライブラリの第1層、第2層、本仮想プロセッサ、C-Threads、および、SunLWPにおいて利用可能である。

測定に用いたプログラムの骨格を、付録D-3に示す。



表 3 - 3 コンテキスト切替えを伴うプロセス間の同期の実行時間の比較

上段：同期操作の実行時間（単位：μ秒）。

下段：繰返し回数（単位：回）。

Table 3-3: Comparisons of execution times of inter-process synchronization with context switches.

Upper: Execution times of synchronization operations. (In microseconds)

Lower: The numbers of iterations.

	SPARCstation 2	Sun3	NeXT	Luna88k	Balance 8000
Microprocess Layer 3/mon(1)	265.8 (10k)	322. (1k)	305. (1k)	66. (10k)	1660. (1k)
Microprocess Layer 3/mon(2)	126.4 (10k)	276. (10k)	246.5 (10k)	48. (10k)	1220. (1k)
Microprocess Layer 3/sem	120.9 (10k)	240. (1k)	184. (1k)	63. (10k)	1700. (100)
Microprocess Layer 2	99.2 (10k)	164. (10k)	154. (10k)	31.4 (100k)	740. (1k)
Microprocess Layer 1	42.1 (10k)	27.1 (10k)	25.45 (100k)	7.3 (100k)	182. (10k)
SunLWP (monitor)	80.2 (10k)	236. (1k)	231. (1k)	— —	— —
C-Thread (mutex)	— —	— —	1146. (1k)	170. (1k)	— —
UNIX (semaphore)	646 (1k)	3200. (1k)	— —	— —	7400. (100)

mon: monitor  
sem: semaphore

実験では、2つのプロセスを用いて、これらの操作を繰り返し行い、実行時間を測定した。実験の結果を、表 3 - 3 に示す。これらの表の上段は、同期操作 1 回当たりの実行時間、下段は、その実行時間を測定した時に行った操作の回数である。ここでは、コンテキスト切り替え 2 回（往復）合わせて 1 回の操作と数える。すなわち、(1) では、セマフォ・プリミティブが 2 つのプロセスで合計 4 回、(2) では、条件変数

の操作が2つのプロセスで合計4回、(3)では、コンテキスト切り替えプリミティブが2つのプロセスで合計2回実行される。

共有メモリ型マルチプロセッサでは、これらの同期プリミティブを発行したとしても、実際には、コンテキスト切替えが起こらない可能性がある。すなわち、2つのプロセスに2つの実プロセッサが割り当てられ、プロセスとプロセッサの対応関係が変化しない可能性がある。しかしながら、利用した性能測定用プログラムにより、これらの同期プリミティブの本質的なオーバーヘッドの測定と実現方式の比較を行うことが可能である。

#### ■ マイクロプロセスの基本的な性能

本実現 Microprocess の各層の機能を比較すると、第1層、第2層の性能がよいことがわかる。コンテキスト切替えにおいては、セマフォ sem とモニタ mon(1), mon(2) は、いずれのアーキテクチャにおいても、同程度のオーバーヘッドになっている。プロセッサの性能と比較して、SPARCstation 2 におけるコンテキスト切替え操作のオーバーヘッドが大きくなっている。この理由は、この操作において手続き呼出しの深度が深く、その結果、レジスタ・ファイルのフラッシュの処理時間が長くなっているからである。

#### ■ 他の軽量プロセスとの比較

表3-3において、本実現 Microprocess Layer 3/mon(2) と C-Threads を比較すると、本実現が約5倍高速に高速になっていることがわかる。この差は、軽量プロセスの実現方式の違いによる。カーネル制御方式の C-Threads では、コンテキスト切替えが必要となる場合、必ずカーネル・コールを発行しなければならないのに対して、本方式のマイクロプロセスでは、その必要がない。

本実現 Microprocess Layer 3/mon(1) と SunLWP を比較すると、本実現の方が効率が悪くなっている。この差は、軽量プロセス実現方式の違いによる。本方式では、共有メモリ型マルチプロセッサ上の並列処理を行うので、細かい単位で共有変数の排他制御を行っている。

#### ■ UNIXのプロセスとの比較

表3-3において、本実現 Microprocess Layer 3/sem と UNIX semaphore を比較すると、本実現の方が5倍から13倍高速になっている。

表 3 - 4 仮想プロセッサの生成・消滅の性能  
(単位 : m 秒)

Table 3-4: Performance of virtual processor creation.  
(In milliseconds)

	Luna88k
ReSC Process	15.
Virtual Processor	7.7
Microprocess (Layer 2)	0.19
UNIX Process	16.
Mach C-Thread	10.5

表 3 - 5 仮想プロセッサのコンテキスト切替えの性能  
(単位 : μ 秒)

Table 3-5: Performance of virtual processor creation.  
(In microseconds)

	Luna88k
ReSC Process	121.
Virtual Processor	68.
Microprocess (Layer 2)	31.
UNIX Process	
Mach C-Thread	170.

### 3. 1 2. 6 仮想プロセッサの性能

実現したシステムの性能を、表 3 - 4 と表 3 - 5 に示す。実現システムのプロセスと仮想プロセッサ (Virtual Processor) は、開発システムとして用いた Mach2.5 システムのプロセス (UNIX prrocess)、および、Mach C-Threads と同程度の性能が得られることが確認された。

利用者レベルの軽量プロセスであるマイクロプロセスと比較すると、仮想プロセッサは、生成において 4 0 倍、コンテキスト切替えにおいて 2 倍程、実行時間が大きくなっている。これは、カーネル・コールに伴うトラップ処理とパラメタ・チェックのオーバーヘッドによる。

### 3. 1 2. 7 応用固有のスケジューラの性能

3. 8. 3 項と 3. 9. 2 項において述べたデータベースの並列処理システム S M A S H の応用固有のスケジューラを記述した。さらに、そのスケジューラの効果

を調べる実験を行った。

実験に用いた問い合わせの構造を、図3-20に示す。このように、8個の軽量プロセスより構成されている。

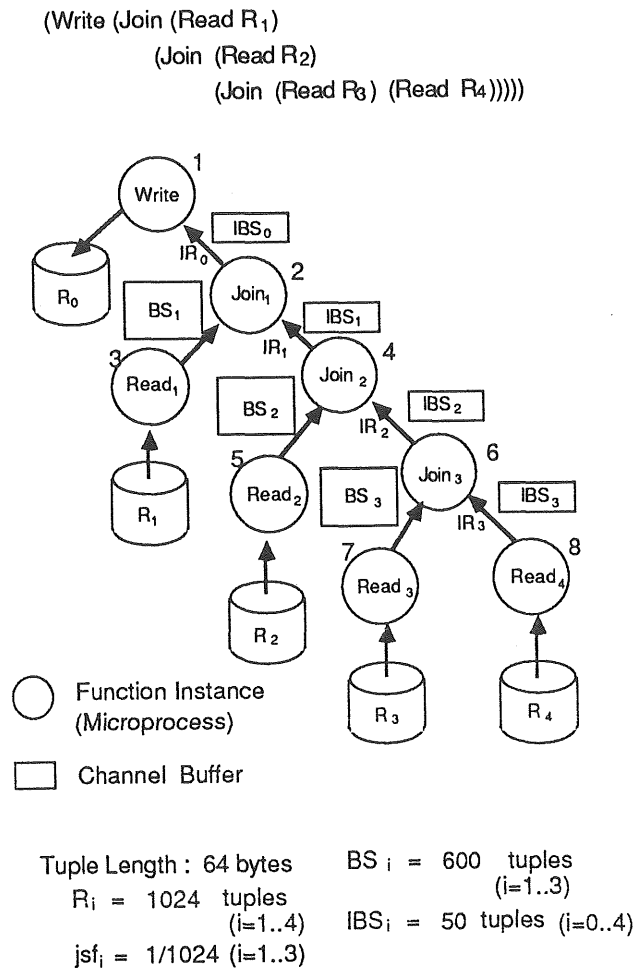


図3-20 実験に用いた関係データベースへの問い合わせの構造

Figure 3-20: The structure of a query for relational databases.

表3-6に実験結果を示す。表の結果より、優先順位を用いないスケジューラ(1)と比較して、優先順位を用いるスケジューラである(2)、または、(3)の性能がよいことがわかる。特に、利用可能なプロセッサの数が少ない場合において、優先順位を用いるスケジューリングの効果が大きいことがわかる。(この実験では、他にプロセスが存在しないので、仮想プロセッサの数と利用可能な実プロセッサの数は、一致している。)これは、動的に優先順位を決定することにより、無駄な計算を避けることができたからである。

以上のことから、本方式において、今回利用した並列応用プログラムについて、応用固有のスケジューラが容易に記述できること、および、そのスケジューラにより効率が改善されることが示された。

表 3 - 6 データベースの並列処理の実行結果 (単位 : 秒)

Table 3-6: Execution times of a parallel application program  
(database processing). (in seconds)

# of active VP	Balance 8000			Luna88k		
	Scheduler 1	Scheduler 2	Scheduler 3	Scheduler 1	Scheduler 2	Scheduler 3
1	40	38	38	4.5	3.7	3.7
2	22	21	26	3.2	2.9	2.7
3	18	16	19	3.0	3.0	2.9

### 3. 1 3 関連した研究

#### 3. 1 3. 1 他の軽量プロセス実現方式との比較

他の実現方式との比較について、3. 4. 1項では、その概略を述べた。

3. 5. 2項では、セマフォ同期プリミティブの実現を通じて、比較を行った。

3. 1 2. 4項、3. 1 2. 5項、3. 1 2. 6項では、性能面の比較を行った。この項では、簡単にまとめを行う。また、これらの項で触れなかった面についても、比較を行う。

#### ■カーネル制御方式との比較

カーネル制御方式と比較して、本方式の利点は、効率的な実現が可能である点にある。この理由は、本方式では、軽量プロセスをサブルーチンと同等に扱い、保護の機能を設けていないことによる。保護とは、軽量プロセス操作を実現している手続き、または、カーネル・コールの入口において、パラメタのチェックを行うことで、不正な操作から利用者プロセスが破壊されないようにすることである。カーネル制御方式では、パラメタのチェックを省略することができないため、結果として同一プロセス内において保護が行われていることに相当する。また、本方式の方が、軽量プロセスの数に関する制約が緩いという利点がある。本方式では、メモリ容量の許す限り多くの軽量プロセスを生成することができる。

カーネル制御方式の利点は、横取りを行うスケジューリング、および、実時間性を保証するスケジューリングの実現が容易である点にある。また、複数のプロセスから構成される応用プログラムを記述する場合も、カーネル制御方式が有利である。それは、カーネルが別のプロセスの中に直接、軽量プロセスを生成する機能を提供することができるからである。本方式では、プロセス間通信を行い、間接的に軽量プロセスを生成することになる。

複数のプロセスを用いて処理を行う場合、カーネル制御方式の方が適している。これは、カーネルにおいて軽量プロセスの優先順位を一元的に管理することが可能であるからである。本方式の場合、異なるプロセスに属する軽量プロセスの間に論理的な関係を記述することができない。この点は、1つのプロセス（資源割当てと保護の単位）を用いて並列処理を行う場合には、なんの問題にもならない。しかしながら、クライアント・サーバ間の遠隔手続き呼出しの実行のように、複数のプロセス間の論理的な処理単位が相互作用を行う時には、カーネル制御方式の方が有利である。

本方式における仮想プロセッサは、カーネル制御方式の軽量プロセスと類似している。相違点は、本方式では、軽量プロセスが利用者レベルにおいて実現されていることを前提にして、機能が設計されていることである。たとえば、カーネル制御方式においては、各軽量プロセスに固有領域を設ける必要はない。本方式の仮想プロセッサ

では、利用者レベルにおいて軽量プロセスを実現するために、仮想プロセッサに固有領域を設ける必要がある。また、カーネル制御方式の軽量プロセスでは、スケジューリングにおいて優先順位が重要な役割を果たすのに対して、本方式の仮想プロセッサでは、優先順位概念が存在しない。

プロセッサのアーキテクチャによっては、カーネルにおいて制御を行った方が効率がよい場合がある。たとえばSPARCプロセッサ [93] [71] では、レジスタ・ファイルのフラッシュを特権モードで行う必要があるため、コンテキスト切替えにおいて必ずカーネル・コールが必要となる。この場合、カーネルにおいてコンテキスト切替えを制御する方が効率が良いことが予想される。

### ■ コルーチン方式との比較

コルーチン方式と比較して、本方式の利点は、共有メモリ型マルチプロセッサにおいてCPU処理の並列実行が可能である点にある。コルーチン方式の利点は、実行時のオーバーヘッドが本方式よりもさらに小さい点にある。この理由は、コルーチン方式では、動作が完全に逐次的であるため、軽量プロセス間の共有変数の相互排除を行う必要がないからである。

本方式において、相互排除を行っている部分を削除することにより、完全にコルーチンに還元することができる。すなわち、並列処理を行う必要がない場合には、本方式においても、コルーチン方式と同等の性能を得ることが可能である。

コルーチン方式の利点の1つとして、プログラム実行の再現性がある点があげられる。本方式では、3. 8. 2項で述べたように、この性質を利用して、効率的にマイクロプロセス間の同期・通信プリミティブの開発を行うことができる。

### 3. 13. 2 他の仮想プロセッサ方式を用いる方式との比較

カーネルにおいて仮想プロセッサを提供し、利用者レベルの軽量プロセスを実行する方式は、文献 [88] [82] [81] において提案した本方式の他にも提案された [3] [23] [30] [53]。これらの方式と比較して、本方式の特徴を、以下にまとめる。

(1) 利用者レベルの軽量プロセスが、層構造を持つライブラリにより提供されている。これにより、各並列プログラムのプログラマは、それぞれの要求に応じてシステムの機能を利用することができる。

(2) 本方式では、コルーチンの性質を利用して、軽量プロセス間の同期・通信プリミティブを構築することができる。

(3) 本方式では、並列応用プログラムとほとんど同じ技術により、仮想プロセッサを実現するカーネルを構成している。これにより、カーネル上で動作する並列応用プログラムとしてカーネルの一部を実行し、動作の確認とデバッグを行うことができる。

(4) 仮想プロセッサと実プロセッサに固有領域がある。これにより、それぞれ並

列応用プログラムとカーネルのプログラムが書きやすくなっている。この固有領域固有領域とカーネル内軽量プロセス・スケジューラによる仮想プロセッサとプロセスのスケジューラを用いる方式は、異なる外部仕様を持つ仮想プロセッサ、あるいは、カーネル・レベルの軽量プロセスを提供するカーネルの実現においても利用することが可能である。たとえば、スケジューラ活動体 (scheduler activations) [3] や Mach の Threads [2] と同じ外部仕様をもつものを 3.10 節で述べた固有領域や 2 段階のレディ・キューを持つスケジューラを用いて実現することも可能である。ただしこの場合、カーネルと利用者プロセスの構造の対称性は、失われる。

(5) 本方式では、カーネルと利用者プログラムが分散的に動作する。すなわち、カーネルと利用者間に共有領域がなく、かつ、カーネルと利用者間の相互作用は、小数のカーネル・コールに限られている。これにより、カーネルと利用者プログラムの双方が書きやすくなっている。しかしながら、共有領域を用いることにより、効率的にカーネルと利用者間の情報交換を行うことが可能になる局面が存在する。たとえば、カーネルから利用者に対して提供する情報としては、文献 [30] において提案されている横取りの情報や、3.5.3 項において述べた時刻の提供などがあげられる。利用者からカーネルへの情報としては、文献 [23] のように、現在 (実プロセッサが割り当てられている) 仮想プロセッサが実行している軽量プロセスの優先順位があげられる。

### 3.13.3 スピンロックとプロセッサの横取りの問題

多重プログラミング・システムにおいて軽量プロセスを使って並列処理を行う場合、スピンロックとプロセッサの横取りの問題を解決する必要がある。この問題とは、スピンロックにより実現された際どい部分 (critical section) を実行中に実プロセッサが横取りされることにより生じる問題である。実プロセッサが割り当てられている仮想プロセッサが際どい部分に入り、スピンロックを保持したとする。この時、プロセス・スケジューラ等により、その仮想プロセッサから実プロセッサが横取りされることがある。この時、他の仮想プロセッサがそのスピンロックを保持しようと試みると、大量の CPU 時間がスピンのために浪費されることになる。

この方式への対応方法は、大きく次の 3 つに分類される。

#### (1) 防止 (prevention)

1 量子時間には、1 つの並列応用プログラムにそのプロセッサ (サイト) に存在する全ての実プロセッサを割り当てる。1 つのプロセスでは、実プロセッサが同時に横取りされるため、上記の問題が防止される。

#### (2) 回避 (avoidance)

利用者プログラムに対して、カーネルが、カーネル・レベルで横取りが行われたこ



とをソフトウェア割込みや上向き呼出し (upcall) により告知する。利用者プログラムは、カーネルからの情報を元に、上記の問題を回避するようにスケジューリングを行う。たとえば、スケジューラ活動体 [3] では、カーネル・レベルにおいて、スケジューラ活動体 (一種の仮想プロセッサ) から実プロセッサが横取りされた場合、(最後の1個の時を除いて) 新たにスケジューラ活動体を生成し、他の実行中のスケジューラ活動体から実プロセッサを横取りし、利用者に報告する。この報告は、カーネルから利用者への1種の上向き呼出しにより行われる。この時、引数として2つのスケジューラ活動体の情報が利用者に渡される。2つのスケジューラ活動体とは、1つは、最初に実プロセッサが横取りされたもの、もう1つは、利用者への報告のために、実プロセッサが横取りされたものである。新たに生成されたスケジューラ活動体では、直前に動いていた2つのスケジューラ活動体の情報を解析し、動作していた利用者レベルの軽量プロセスや、際どい部分 (critical section) の情報を得る。この時、スピロックによる際どい部分を実行していたならば、その部分へ制御を移す。このようにして、スピロックとプロセッサの横取りの問題を回避する。

### (3) 検出 (detection) と解消

スピロックのルーチンの中で、上記の問題が生じたことを検出し、解消を行う。

本方式では、文献 [40] と同様に、(3) より対応する。3.6.2項において述べたように、スピロックを実現するルーチンにおいて、この問題が生じたことを検出する。そして、カーネル・コール `vp_switch()` を発行することにより、他の仮想プロセッサへ制御を移し、この問題の解消する。

本方式において (3) を採用した理由を、以下に示す。

- ・本方式では、3.4.2項において述べたように、1つのプロセスの内部は、横取りを行わないスケジューリングを行う。従って、この問題が生じる可能性が小さい。
- ・3.12.6項において示したように、この問題が生じた時に復旧するための代償が小さい。この問題が生じた時に起きる仮想プロセッサ間のコンテキスト切替えは、ソフトウェア割込みやスケジューラ活動体の生成と比較して、オーバーヘッドが小さい。

### 3.13.4 共有 fork との比較

3.6.3項で述べたアドレス空間の構造は、D y n i x システム [6] における fork システム・コール、および、S p r i t e システム [60] における共有 fork (shared fork) の結果と類似している。相違点は、D y n i x や S p r i t e の場合、新たにプロセスが作られる点にある。すなわち、元のプロセスと新たに生成されたプロセスは、異なるプロセス識別子を持ち、それらが協調して動作していることは、外部から観測している限り分らない。本方式では、プロセスと仮想プロセッサは、異

なるレベルにあり、任意のプロセスを、内部の仮想プロセッサの数にかかわらず一様に操作することが可能である。たとえば、プロセスが終了すると、内部のすべての仮想プロセッサも自動的に消去される。これに対して、D y n i x や S p r i t e では、個々のプロセスが別々に終了する。

### 3. 1 3. 5 軽量プロセスの集合によるカーネルの構築

分散型オペレーティング・システムを軽量プロセスの集合として実現する方法は、文献 [ 9 7 ] において既に提案されている。その方法では、相互排除アクセスされる資源の管理単位に軽量プロセスを配置する。軽量プロセス間のデータのやりとりは、軽量プロセス間の共有メモリではなく、プロセス間通信で行う。このため、オペレーティング・システムのカーネル自身をネットワーク上の別のプロセッサに配置することが可能となっている。

この構成法と 3. 1 0 節で述べた本構成法の最大の違いは、本構成法では、並列処理の単位に軽量プロセスを配置する点にある。また、本構成法では、軽量プロセスが同一アドレス空間に存在することを利用して、軽量プロセス間のデータのやり取りを軽量プロセス間の共有メモリを用いて行う。

### 3. 1 4 まとめ

この章では、マイクロプロセスと仮想プロセッサの概念に基づく軽量プロセスの実現方式について述べた。本方式では、利用者レベルにおいて軽量プロセスを実現することにより、高い性能を得ることができる。さらに、利用者プロセス内の軽量プロセスのスケジューラにより、応用固有のスケジューリングが実現される。本システムでは、層構造をもつライブラリにより軽量プロセスの実現を支援する。データベースの並列処理システムの構築を例として、そのライブラリを用いた応用固有の軽量プロセス間の同期・通信プリミティブの開発や、応用固有のスケジューラの開発の方法について述べた。

次に、仮想プロセッサを提供するカーネルの構成法について述べた。本構成法の特徴は、利用者プロセスとカーネルが同一の構造を持っている点にある。両者とも軽量プロセスから構成される。そして、利用者レベルの軽量プロセスは、仮想プロセッサにより、カーネル内の軽量プロセスは、実プロセッサにより実行される。各仮想プロセッサと各実プロセッサは、固有のメモリ領域を持ち、そこにそれぞれの識別子や固有のデータを保持する。仮想プロセッサは、カーネルの内部において並列処理の単位であり、カーネル内の軽量プロセスとして実現される。プロセスは、仮想プロセッサを実現している軽量プロセスの集合として実現される。利用者レベルの抽象であるプロセスと仮想プロセッサのスケジューリングは、カーネル内の軽量プロセスのスケジューラにより行われる。そのスケジューラは、プロセス・レベルと仮想プロセッサ・

レベルの2段階のレディ・キューから構成される。

提案した軽量プロセス実現方式、および、カーネル構成法に従って、軽量プロセスを実現するライブラリと仮想プロセッサを提供するカーネルの実現を行った。そして、その基本的な性能を測定した。本方式の軽量プロセスと、カーネル制御方式、および、コルーチン方式の軽量プロセスの基本的な性能を比較する実験を行った。さらに、1つの並列応用プログラムについて応用固有のスケジューラの記述とその性能を調べる実験を行った。これらの実現と実験を通じて、本方式の有効性を示した。

## 第4章 マッピング・コントローラ

この論文において、マッピングとは、ネットワークに結合された複数のプロセッサに対して、プロセスやデータを割り当てることを意味するものとする。マッピング・コントローラとは、1つの応用プログラム中にある、マッピングの方針を決定するモジュールである。この章では、マッピング・コントローラの基本概念について整理を行い、R e S Cシステムにおけるマッピング・コントローラに対する方針と支援について述べる。次に、並列シェルの機能とそのプロトタイプを用いた実験について述べる。並列シェルは、内部にマッピング・コントローラを持たない応用プログラムに代わり、マッピングの最適化を行うものである。

### 4.1 マッピング・コントローラ概念

ネットワーク上に分散した複数のサイト（プロセッサ）を用いて並列処理や分散処理を行う応用プログラムにとって、そのプログラムのプロセスやデータのサイトへの割当ては、非常に重要である。このようなプログラムの中にあり、割当て方針を決定するモジュールをマッピング・コントローラ（mapping controller）、あるいは、分散オブティマイザ（distribution optimizer）とよぶことにする [79] [82] [86]。

マッピング・コントローラは、システムの資源の利用状況を観察し、プロセスやデータの最適な割当てを行う。この時、各プロセスの重さやプロセス間の通信量の予測を元に、プロセスやデータ・ファイルの割当ての最適化を行う。たとえば、CPU処理が中心のプログラムのマッピング・コントローラは、CPUの能力が高く負荷が軽いサイトを探し、プロセスを生成する。大きなファイルを検索する応用プログラムのマッピング・コントローラは、ネットワーク通信のオーバーヘッドを削減するために、ファイルが存在するサイトにプロセスを生成する。

マッピング・コントローラは、マッピングの方針（policy）を決定するものである。具体的なデータ・ファイルやプロセスの生成という仕組み（mechanism）は、ファイル・サーバやプロセス・サーバにより提供される。

### 4.2 マッピング・コントローラに対するR e S Cシステムの方針

マッピング・コントローラに対するR e S Cシステムの方針を、以下にまとめる。

- (1) ReSCシステムのカーネルは、マッピング・コントローラを含まない。
- (2) 各並列／分散応用プログラムは、固有のマッピング・コントローラを備える。
- (3) システムは、各応用プログラムのマッピング・コントローラに対してシステム構成や資源の利用状況などの情報を積極的に提供する。
- (4) マッピング・コントローラを含まない応用プログラムに対して、システムは、並列シェルを提供し、ネットワーク上に分散した資源の有効利用を促進する。

2. 4節、および、3. 4節において述べたように、カーネルは、1つのサイト内における応用プログラム間の公平なCPU資源のスケジューリングを行う。しかしながら、サイトへのマッピングに関して、カーネルは、仕組み (mechanism) を提供するファイル・サーバやプロセス・サーバを含んでいるが、方針 (policy) を決定するマッピング・コントローラを含んでいない。すなわち、マッピングに関して方針と仕組みの分離 (policy/mechanism separation) が実現されている [108]。

上記の方針により、以下のような利点が生じる。

- (1) 各応用プログラムが、その応用プログラム固有の性質を利用するようなマッピング・コントローラを構築することができる。
- (2) マッピングを制御する時に、システム定義の汎用の記述ではなく、応用プログラム固有の記述を用いることができる。
- (3) マッピング・コントローラを持っていない応用プログラムに対しても、並列シェルによりネットワーク上の資源の有効利用が実現される。

上記(1)と(2)は、マッピングに関して、システムと並列／分散プログラム間の調和が実現されることを意味する。すなわち、システムは、並列／分散応用プログラムと競合することなく、そのような応用プログラムに対して、高信頼性や高速化のためにマッピングを制御する仕組みを提供している。また、上記(3)は、マッピングに関して、逐次応用プログラムと並列／分散応用プログラムの間の調和が実現されていることを意味する。すなわち、マッピング・コントローラを持たない逐次応用プログラムは、並列シェルを利用して、システムを集中型システムと同等に扱うことができる。同時に、並列／分散応用プログラムは、カーネルを直接利用することで、並列シェルに妨害されることなく、システムを分散型システム利用することができる。

#### 4. 3 応用固有のマッピング・コントローラの支援

単一プログラミング・システムでは、各並列／分散応用プログラムは、全ての資源の利用状況を常に把握することができる。マッピング・コントローラは、この情報を利用して、最も効率的な実行が可能なように、プロセスの配置やデータの配置の最適化を行う。一方、多重プログラミング・システムにおいては、オペレーティング・システムがそのような情報を提供しなければならない。それは、システムの利用状況が刻々と変化するからである。

ReSCシステムは、以下の情報を応用固有のマッピング・コントローラへ提供す

る。

- (1) 各サイトのプロセッサの型、数、および、速度
- (2) 各サイトの実メモリ容量
- (3) ネットワークの形態 (topology)
- (4) 各サイトの稼働状況 (動作中かどうか) とプロセッサの負荷
- (5) 各サイトのメモリの利用状況
- (6) ネットワークのトラヒック
- (7) 資源の位置

項目 (1) から (3) までは、静的な情報である。これらは、公のファイルに格納しておく。項目 (4) から (6) は、システムの動的な情報である。これらは、カーネル・サーバにより提供される。項目 (7) は、2. 6 節で述べたオブジェクト識別子のサイト識別子を通じて提供される。

#### 4. 4 並列シェル

並列シェル (parallel shell) は、協調して動作する複数の逐次プログラム (コマンド) のためのコマンド・インタプリタである [79] [82] [86] [36]。並列シェルは、R e S C システムの構成要素であると同時に、カーネル上で動作する 1 つの並列応用プログラムでもある。並列シェルは、内部にマッピング・コントローラを持たない (逐次) 応用プログラムに代わり、マッピング・コントローラとしての役割を果たす。たとえば、その応用プログラムをどのサイトで実行するかを決定したり、出力ファイルを作成するサイトを決定する。

4. 3 節で述べた情報に加えて、並列シェルは、次の様な情報を利用してマッピングの最適化を行う。

- (1) 各コマンドの負荷の予測
- (2) それらの間の通信量の予測
- (3) 結合関係

これらは、並列シェルにとって、応用固有の情報に相当する。すなわち、他の並列／分散応用プログラムには存在しない、並列シェルに特有の概念である。

##### 4. 4. 1 並列シェルの構文と機能

並列シェルの構文と意味 (機能) は、UNIX のシェルのものを拡張したものである。UNIX のシェルから次のような機能と構文を引き継いでいる。

- (1) パイプ (|)
- (2) 入出力の切替え (>, <)
- (3) コマンドの並列実行 (&)

これらの加えて、次のような拡張を行った。

- (1) on 文
- (2) siteof 組込関数
- (3) processor\_allocate 組込関数

on 文は、サイトのリストを受け取り、そのサイトでコマンドを実行する。siteof組込関数は、そのオブジェクトの位置を返す。processor\_allocate 組込関数は、高速で負荷が小さいサイトの識別子を返す。

以下に、on() と siteof() を使った例を示す。このようなシェル・スクリプトが与えられると、並列シェルは、ストリームにより結合された3つのコマンドを起動する(図4-1(a))。

```
on(siteof(infile)) com-1 < infile | ¥
com-2 | ¥
on(siteof(outfile)) com-3 > outfile
```

ここで、“infile” と “outfile” は、そのストリームの入力となるファイルと出力となるファイルの名前である。この例では、入力ファイルが存在するサイト siteof(input) において com-1 を実行し、出力ファイルが存在するサイト siteof(output) において com-3 を実行する。com-2 を実行するサイトは、記述されていない。このような場合、並列シェルにより決定される。

一方、図4-1(b)は、簡単な負荷分散を行うマッピング・コントローラによるコマンドのマッピングの例である。このマッピング・コントローラは、2つのサイトの負荷均衡を試みている。しかしながら、コマンド間の結合関係の情報を利用していないため、結果として、ネットワーク通信量が増加している。

以下に、processor\_allocate() 組込関数の例を示す。

```
on(processor_allocate(1,1,2)) com-1 & com-2 & com-3
```

この例において、processor\_allocate() 関数は、負荷の予測を受け取り、軽いサイトのリストを返す。そして、on() 文により、それらのサイト上で3つのコマンドを並列に実行している。

ここで述べた on文、siteof組込関数等は、並列シェルに固有のマッピングの記述方法である。

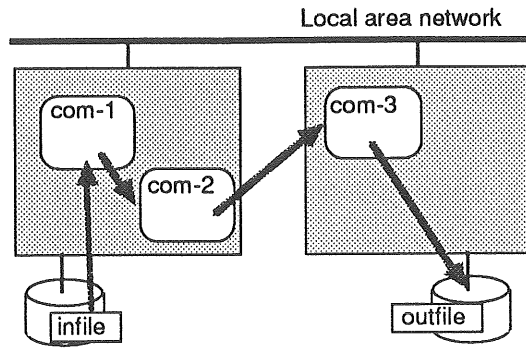


図 4 - 1 ( a ) 並列シェルによるもの  
 Figure 4-1(a): By the parallel shell.

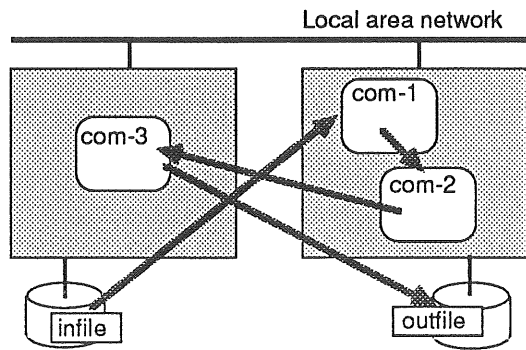


図 4 - 1 ( b ) 簡単な負荷分散を行うマッピング・コントローラによるもの  
 Figure 4-1(a): By a simple mapping controller.

図 4 - 1 ストリームで結合されたコマンドのマッピング  
 Figure 4-1: Mapping commands connected with streams.

#### 4. 4. 2 プロトタイプの実現

並列シェルのプロトタイプを、イーサネットにより結合されたワークステーション上に実現した。実現においては、SunRPC [69] を利用して、2. 7 節において述べた ReSC システムの遠隔手続き呼出しのエミュレーションを行った。ネットワーク上の各サイト（ワークステーション）には、ReSC のプロセス・サーバのエミュレータを配置する。並列シェルは、遠隔手続き呼出しを用いて、各サイト上のプロセス・サーバに対してプロセス生成要求を送る。要求を受け付けたプロセス・サーバは、UNIX の fork システム・コール、および、execve システム・コールを用いて、プロセスを生成する。パイプの実現には、TCP/IP [101] のストリームを利用した。



### 4.4.3 実験

実現した並列シェルを用いて、コマンド間の結合関係の情報を利用することによる処理効率の違いを調べる実験を行った。実験に用いたシェル・スクリプト (shell script, 並列シェル・インタプリタのプログラム) の性質を、表4-1に示す。ここでは、両スクリプトともストリームにより結合された8個のコマンドを含んでおり、それらのコマンドは、全体でパイプラインを構成する。パイプラインの各コマンドは、1ブロック (1kバイト) のデータを入力し、そのデータに関して処理を行い、同じ大きさのデータを出力する。パイプラインの先頭のコマンドの入力は、読み出すと0で埋められたデータを返す特殊なファイル (/dev/zero) である。パイプラインの末尾のコマンドの出力は、書き込んだデータを吸い込む特殊なファイル (/dev/null) である。

表4-1 実験に用いたスクリプトの性質

Table 4-1: The behavior of experimental scripts.

	Script 1	Script 2
Load (m sec/KB)	2	200
IPC (KB)	200	2

スクリプト1では、各コマンドの入出力データに対する計算深度が浅い。すなわち、各コマンドは、入力したデータに対して、単純な処理を行う。プロセス間の通信量も大きい。そして、パイプラインの定常状態が長い。スクリプト2は、スクリプト1とは逆に、計算深度が深く、プロセス間通信量が小さい。パイプラインの定常状態は、存在しない。

実験に用いたマッピングを以下に示す。

```
s0, s1=processor_allocate(1, 1)
```

Sequential:

```
on($s0) source | p | p | p | p | p | p | sink
```

Parallel(a):

```
on($s0) (source | p | p | p) | Y
on($s1) (p | p | p | sink)
```

Parallel(b):

```
on($s0) source | on($s1) p | on($s0) p | on($s1) p | Y
```

on(\$s0) p | on(\$s1) p | on(\$s0) p | on(\$s1) sink

このスクリプトにおいて、括弧は、UNIXのシェルとは異なり、結合性を変えることを意味し、新しいシェルを作ることの意味しない。

表4-2に、実験結果を示す。このように、両スクリプトのいずれのマッピングにおいても、並列処理により実行時間が短縮されていることがわかる。スクリプト1では、マッピング Parallel(a)の方が実行時間が短い。これは、ネットワーク通信のオーバーヘッドが削減されたことによる。対照的に、スクリプト2では、マッピング Parallel(b)の方が実行時間が短い。この理由は、このスクリプトではパイプラインの定常状態がなく、実行の途中においても、Parallel(b)では、2つのサイトの負荷が均衡しているからである。Parallel(a)では、実行の前半では、サイト \$s0 に、後半では、サイト \$s1 に負荷が集中している。

このように、スクリプトの性質により、最適なマッピングが異なる。このことは、スクリプトの中のコマンド間の結合関係の情報を利用することにより、資源割当ての最適化が可能であることを示している。

表4-2 スクリプトの実行時間

Table 4-1: Execution times of the scripts.

mapping	Script 1	Script 2
Sequential	6.2	5.5
Parallel (a)	3.2	3.9
Parallel (b)	4.0	3.6

## 4.5 関連した研究

### 4.5.1 従来の分散型オペレーティング・システムとの比較

従来の分散型オペレーティング・システム [98] [22] [43] では、システムがマッピング・コントローラを含んでいる。すなわち、システムがマッピング・コントローラとして、マッピングの方針 (policy) を決定することを意味する。このため、並列/分散応用プログラムを走らせると、システムのマッピング・コントローラと各応用プログラムのマッピング・コントローラが競合することになる。これに対して、本 R e S C システムでは、カーネルのレベルではマッピング・コントローラを含んでいない。このため、従来の分散型オペレーティング・システムにみられたような競合が解消される。

A m o e b a システムは、従来の定義と目標そった典型的な分散型オペレーティン

グ・システムである [ 5 4 ] 。 ( 本システムにおける分散型オペレーティング・システムの定義と目標については、1. 5 節において述べた。 ) Amoeba システムは、2. 2. 1 項で述べたサーバ・モデルにより構築されているため、ネットワークの位置を利用する状況が少ない。すなわち、ファイル・サーバやウインドウ・サーバのプロセスを除き、応用プログラムは、全てプロセッサ・プールとよばれるネットワークで結合された計算サーバ群の 1 つで実行される。システムは、プロセッサ・プールの利用状況を常に監視しており、応用プログラムからのプロセッサ割当て要求に対して、負荷が小さいプロセッサを割り当てる。すなわち、システムがマッピング・コントローラを含んでいる。

Amoeba では、プロセッサ・プールとそれを管理するオペレーティング・システムの機能により、従来の逐次的なプログラムを複数個、並列に実行し、システム全体の処理能力を改善することを試みる。しかしながら、応用プログラムのレベルでマッピング・コントローラを支援する機能が存在しない。これに対して、本システムでは、並列／分散応用プログラムのマッピング・コントローラに対して、資源の利用状況や位置の情報を積極的に提供する。

Amoeba システムにおけるマッピングの機能は、本システムでは、並列シェルのレベルにおいて提供される。ゆえに、本システムの並列シェルとカーネルは、従来の定義における分散型オペレーティング・システムを構成するための方法と見ることも可能である。

Amoeba システムでは、プロセス生成のレベルにおいて、プロセスを単位としてマッピングの最適化を行う。これに対して、本システムでは、並列シェルのレベルにおいて、複数のコマンドのプロセスの結合関係を利用して、マッピングの最適化を行う。マッピングの最適化を行う場合、シェルのレベルで行う方が、プロセス生成のレベルで行うよりも、より簡単で効果的な最適化が可能になる。それは、シェルのレベルでは、コマンドの入出力ファイルの位置や、コマンド間の結合関係を得ることが容易であるからである。プロセス生成のレベルでは、そのプロセスがどのようなファイルを読み書きするのか、あるいは、他のプロセスとどのようなプロセス間通信を行うかを調べることは、困難である。

#### 4. 5. 2 負荷均衡の方針の研究

4. 4 節で述べた並列シェルは、主に負荷均衡の仕組みを提供するものであり、それを用いてさまざまな方針のマッピングを記述することが可能になっている。マッピングの方針としては、様々な研究が行われている [ 2 2 ] [ 4 3 ] 。並列シェルには、これらの研究成果を取り入れることが可能である。

文献 [ 3 6 ] では、独立に実行可能なコマンド群とパイプラインを形成しているコマンド群について、簡単なマッピングの方針を提案した。前者の特徴は、受け付けたコマンドを全て一度に実行するのではなく、最大のスループットが得られるように各サイトの多重度を調整しながらプロセスを割り当てる点にある。そして、プロセスが

終了したサイトに、順次プロセスを割り当てることにより、各プロセスの実行時間の違いや動的な負荷の変動への適応が実現される。後者については、4.4.3項の実験で用いたような2種類のマッピングの方針（サイト間通信を減らすものとパイプライン処理を考慮した負荷均衡を図るもの）を提案した。

#### 4.5.3 負荷均衡の仕組みの研究

4.4節において述べた並列シェルとプロセス・サーバは、本ReSCシステムの一部として設計されたものである。そのプロトタイプは、SunRPCとTCP/IPによる通信機能があるワークステーションにおいて利用可能である。したがって、本システムの並列シェルとプロセス・サーバを、そのような環境における負荷均衡の仕組みを提供するものとしてとらえることが可能である。

他の負荷均衡の仕組みの研究[22][19]と比較した場合、本並列シェルの特徴は、ストリームにより結合されたコマンド群の並列実行が可能である点にある。他の多くの研究では、主に独立して動作する（プロセス間通信を行わない）複数のプロセスが対象となっている。

#### 4.6 まとめ

この章では、マッピング・コントローラ概念を整理し、ReSCシステムにおけるマッピング・コントローラに対する方針と支援について述べた。ReSCシステムでは、カーネルがマッピング・コントローラを含んでいない。カーネル上で動作する、ネットワーク上の複数のプロセッサ（サイト）を利用する各並列／分散応用プログラムは、内部に固有のマッピング・コントローラを持つことになる。これにより、システムと並列／分散応用プログラムの間、マッピングの方針に関する競合が解消され、調和が図られる。

並列シェルは、マッピング・コントローラを持たない応用プログラムに代わって、マッピングの最適化を行う。シェル・スクリプトのレベルでマッピングの最適化を行うので、プロセス生成のレベルでマッピングの最適化を行う方法と比較して、より高度な最適化が容易に達成される。

## 第5章 オブジェクトの堆積

この章では、オブジェクトの堆積というモデル化の手法について述べる。ReSCシステムにおいて、オブジェクトの堆積は、外部サーバが提供する高度な機能を統合するために用いられる。そして、統合された機能は、主に逐次応用プログラムにより利用される。

この章では、はじめに、オブジェクトの堆積の目標と背景となる技術について述べる。次にオブジェクトの堆積の基本概念を、例と比較を用いて説明する。そして、分散型オペレーティング・システムにおけるオブジェクトの堆積の利用について述べる。

### 5.1 目標と背景

オブジェクトの堆積は、object-based システムにおいて、複数のサーバの持つ機能を統合して利用するためのモデル化の手法である [78] [80] [83]

[85]。オブジェクトの堆積は、単純なモデルであるが、非常に強力である。本ReSCシステム以外のシステムにおいても、広く利用することが可能である。

2章で述べたように、ReSCでは、カーネルが、分散核としてプロセス間通信機能を提供する。同時に、カーネルは、ファイル・サービスやプロセス・サービスなど基本的な機能を提供する。これらの機能は、カーネルで実現されるため、オーバーヘッドが小さく、効率を重視する並列／分散応用プログラムによる利用に適している。しかしながら、逐次応用プログラムの高度な要求を全て満足することはできない。たとえば、ファイルのキャッシングやファイルの複製 (replication)、オブジェクトの移動 (migration) といった機能を提供していない。ReSCでは、このような高度な機能は、カーネル外のサーバにおいて提供される。オブジェクトの堆積は、この高度な機能を実現するモデル化の手法として開発されたものである。

オブジェクトの堆積により、逐次応用プログラムと並列／分散応用プログラムの間の調和が実現される。すなわち、逐次プログラムは、オブジェクトの堆積により統合された、高度な外部サーバの機能を利用することができる。同時に、並列／分散応用プログラムは、効率的なカーネル・サーバを直接利用することができる。また、システムの一部である外部サーバは、並列／分散応用プログラムとしてカーネルの機能を利用する。これにより、システムと他の並列／分散応用プログラムの間の調和が図られる。

オブジェクトの堆積の背景となった技術を、以下にまとめる。

### (1) object-based 分散型オペレーティング・システム

最近の分散型オペレーティング・システムの多くは、object-basedシステムとして構成されている [98] [22] [43]。オブジェクトの堆積は、このような object-basedシステムにおいて、利用可能である。object-based システムは、サーバ、クライアント、カーネルから構成される。サーバは、オブジェクトを管理し、クライアントへサービスを提供する。主要なサービスとして、ファイル・サービス、プロセス・サービス、ディレクトリ・サービスがあげられる。カーネルは、プロセス間通信と物理的な入出力機能を提供する。オブジェクトの堆積では、システムによって提供される基本的なオブジェクト、オブジェクト識別子、および、クライアント・サーバ間のプロセス間通信機能が必要となる。

一方、分散ファイル・システムに関する研究が盛んに行われている [13] [24] [47] [94]。このような研究の多くは、複製 (replica, コピー) 間の一貫性を保証する効率的なアルゴリズムを与えるものである。オブジェクトの堆積は、そのようなアルゴリズムを与えるものではなく、分散ファイル・システムをカーネル外で実現することを可能するためのモデル化の手法を提供する。さらに、オブジェクトの堆積を使うことにより、汎用の複製ファイル・サーバを構築することができる。従来のファイルの複製の研究では、1種類のファイルに対してのみ複製を作ることができた。これに対して、オブジェクトの堆積では、様々な種類のファイルの複製を作ることが可能になる。

### (2) 機能を組み合わせる仕組み

オペレーティング・システムにおいて、機能を組み合わせる仕組みとして、もっとも成功した例の1つが、UNIXのパイプである。パイプを使うことにより、単純なコマンドを組み合わせ、複雑なコマンドを作ることが可能となる。これに対して、オブジェクトの堆積では、単純なオブジェクトを組み合わせ、複雑なオブジェクトを作ることが可能となる。

### (3) 層化プログラミング

層化プログラミングでは、巨大なシステムを複数の層に分割して開発することで、再利用可能性が高いモジュールが構築される。たとえば、ISOのOSI参照モデルでは、通信プロトコルを7つの層に分割することで、異機種間の通信を可能にしている [16] [92]。これに対して、オブジェクトの堆積では、再利用可能性が高いモジュールをサーバとして開発することが可能となる。たとえば、キャッシング機能と複製機能があるファイル・サーバを構築する場合、キャッシング機能しかないサーバと複製機能しかないサーバを組み合わせることで、両者の機能を併せ持つサーバを構築することができる。

## 5. 2 基本概念

オブジェクトの堆積は、次の4つの基本概念から構成される。

### (基本概念1) オブジェクト

オブジェクトの堆積において、オブジェクトとは、手続きとデータがカプセル化されたものである。オブジェクトのデータの参照と更新は、外部に公開された手続きによってのみ行われる。object-based 分散型オペレーティング・システムでは、オブジェクトは、サーバ・プロセスによって管理と保護が行われる。図5-1において、Server A は、ファイル・サーバであり、ファイル・オブジェクトを管理している。そして、クライアントからの読み込み手続き (read()) と書き込み手続き (write()) を受け付け、実行する。そして、それ以外の不正な操作からオブジェクトを保護している。

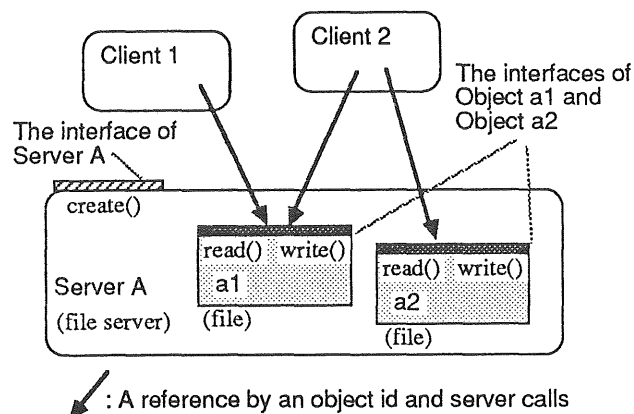


図5-1 オブジェクト、サーバ、クライアントの関係

Figure 5-1: Objects, their server and clients.

オブジェクトの堆積では、オブジェクトは、オブジェクト識別子 (object identifier, object id) により常に間接的に参照される。図5-1において、Client 1 は、Object a1 のオブジェクト識別子を保持しており、それを使ってそのファイルに読み込み/書き込み要求のメッセージを送る。そのメッセージは、サーバに送られ、サーバにより解釈される。

オブジェクトに対する手続きをオブジェクト手続き (object procedure)、サーバに対する手続きをサーバ手続き (server procedure) とよぶことにする。これらの手続きは、それぞれオブジェクト指向プログラミング言語におけるインスタンス・メソッド、および、クラス・メソッドに対応する。

あるオブジェクトが受け付けることができる手続きの集合を、オブジェクトのイン

タフェースとよぶ。たとえば、ファイル・オブジェクトは、インタフェースとして読み込み手続きと書き込み手続きを持つ。同様に、サーバが受け付けることができる手続きの集合を、サーバのインタフェースとよぶ。たとえば、ファイル・サーバのインタフェースには、ファイルの生成を行う手続きが含まれている。

### (基本概念2) 積み重ねる、積む、堆積を作る。

オブジェクトa1の上にオブジェクトb1を積み重ねる (to stack b1 on a1) とは、次の2つの条件を満たすことである。

- (1) オブジェクトb1がオブジェクトa1のオブジェクト識別子を保持している。
- (2) オブジェクトb1が自分自信の機能を実現するために、必ずオブジェクトa1を呼び出す。

この時、オブジェクトb1のサーバは、クライアントとしてオブジェクトa1を利用する。

オブジェクトa1の上にオブジェクトb1が堆積可能である (be stackable on) とは、オブジェクトb1の機能を実現するために、十分なインタフェースと機能をオブジェクトa1が持っていることと定義する。(b1が使わないインタフェースと機能をa1が持っていてよい。)

オブジェクトがオブジェクト識別子を含んでいたとしても、必ずしも積み重ねることを意味しない。たとえば、ディレクトリ・オブジェクトは、その中にオブジェクト識別子と文字列の組を保持している。しかしながら、ディレクトリ・オブジェクトは、そのオブジェクト識別子で示されたオブジェクトの上に積まれたオブジェクトではない。なぜならば、そのオブジェクト識別子は、単なるビット列であり、そのディレクトリ・オブジェクト自身の機能を実現するために、そのオブジェクト識別子で示されたオブジェクトへメッセージを送ることがないからである。

### (基本概念3) 堆積をオブジェクトと同じと考える

あるオブジェクトの上に、別のオブジェクトを積む場合、下位層オブジェクトが、システムにより実現された単純なオブジェクトではないかもしれない。すなわち、既にいくつかのオブジェクトが積み重ねられたものの上に、新たにオブジェクトを積むことができる。

別の言葉では、オブジェクトが堆積したのもまたオブジェクトであり、その上に新たにオブジェクトを積み重ねることが可能である。その結果できたのもまたオブジェクトである。今後、堆積によってできたオブジェクトを特に区別して扱う場合、堆積オブジェクト (stack object) とよぶ。単にオブジェクトという場合、単純なオブジェクトと堆積オブジェクトの両方を含むものとする。

あるオブジェクトの下のオブジェクトを下位層のオブジェクト (lower object) とよぶ。あるオブジェクトの上に積まれているオブジェクトを上位層のオブジェクト (upper object) とよぶ。最上位層のオブジェクトを頂上オブジェクト (top object)、最下位層のオブジェクトを基底オブジェクト (bottom object) とよぶ。



基底オブジェクトは、他のオブジェクトに依存していない。そして、他のオブジェクトのオブジェクト識別子を含んでいない。基底オブジェクトは、堆積により作ることができないので、システムにより提供されなければならない。

図5-2において、Object a1 は、Server A によって管理されている基底オブジェクトであり、その上に Server B によって管理されている Object b1 が積まれ、1つの堆積オブジェクトが作られている。これを堆積オブジェクト b1/a1 と表す。Object b1 は、a1 の識別子を保持しており、b1自身の機能を実現するために、a1の機能呼び出す。堆積オブジェクト b1/a1 の上に、さらに Object c1 が積まれ、堆積オブジェクト c1/b1/a1 が形成されている。

堆積オブジェクトの識別には、その頂上オブジェクトのオブジェクト識別子を利用する。たとえば、図5-2の堆積オブジェクト c1/b1/a1 の識別子として、c1 の識別子が使われる。これにより、クライアントは、そのオブジェクトが堆積オブジェクトであるか、単純なオブジェクトであるかにかかわらず、一様に扱うことができる。

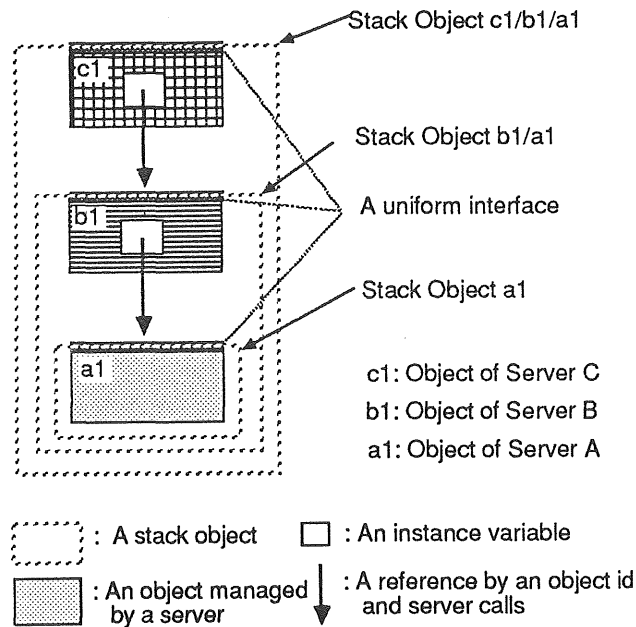


図5-2 サーバに管理されているオブジェクトと堆積オブジェクト

Figure 5-2: Objects managed by servers and stack objects.

#### (基本概念4) 一様なインタフェース

オブジェクトを自由に組み合わせて利用するためには、それらのオブジェクトのインタフェースが一様 (uniform) である必要がある。一様なインタフェースのオブジェクトを利用することにより、積み重ねるオブジェクトの構成と順序を自由に選択することが可能となる。もし、インタフェースが異なれば、あるオブジェクトの上に積むことができるオブジェクトの種類が制約される。また、実行される手続きに着目すると、一様なインタフェースにより、手続きの実行の順序を積み重ねるオブジェクト

の順序として、個々の堆積オブジェクトごとに設定することが可能となる。

一様なインタフェースを持つということは、自分自身のインタフェースと下位層のオブジェクトのインタフェースが同一であることを意味する。たとえば、自分自身がファイルというインタフェースを持つ場合、下位層のオブジェクトのインタフェースもまたファイルである。図5-2において、Object a1, b1, c1 は、一様なインタフェースを持っている。

#### ■堆積可能オブジェクト

自分自身のインタフェースと下位層オブジェクトのインタフェースが一様であるようなオブジェクトを、堆積可能オブジェクト (stackable object) とよぶ。同様に、堆積可能ファイル・オブジェクトとは、自分自身のインタフェースがファイルで、下位層のオブジェクトのインタフェースもファイルであるようなオブジェクトである。図5-2において、Object b1、および、Object c1 が堆積可能である。

堆積可能オブジェクトを管理しているサーバを堆積可能サーバ (stackable server) とよぶ。同様に、堆積可能なファイル・オブジェクトを管理しているサーバを、堆積可能ファイル・サーバとよぶ。図5-2では、Server B、および、Server C が堆積可能サーバである。

### 5. 2. 1 StdFS、ZFS、CFZ

ここでは、以下の3つのファイル・サーバを使って、object-basedシステムにおけるオブジェクトの堆積の例を示す (図5-3)。

#### (1) 標準ファイル・サーバ(StdFS:Standard File Server)

これは、システムが提供する標準的なファイル・サーバである。ReSCシステムでは、カーネルにおいて提供される。標準ファイル・サーバでは、各オブジェクトは、下位層のオブジェクトを必要としない。ファイルのデータは、ディスク・ブロックに格納される。

#### (2) 圧縮ファイル・サーバ(ZFS:Compression File Server)

これは、ファイルの内容を圧縮して保存するようなファイル・サーバである。圧縮ファイル・サーバでは、各オブジェクトは、圧縮したデータを格納するために、下位層のファイル・オブジェクトを1個必要とする。クライアントから書込み要求を受け付けると、サーバは、書込み要求の中のデータを圧縮し (compress)、対応する下位層のオブジェクトに書込み要求を送る。クライアントから読み込み要求を受け付けると、下位層のオブジェクトに読み込み要求を送り、圧縮されたデータを得る。そして、圧縮されたデータを解凍し (uncompress)、その結果をクライアントに返す。

(3) 暗号ファイル・サーバ(CFZ:Encryption File Server)

これは、ファイルの内容を暗号化して保存するようなファイル・サーバである。サーバの構造と動作は、ZFSにおける圧縮と解凍を暗号化と復号化に置き換えたものである。

ZFSでは、普通、各オブジェクトは、下位層のオブジェクトとしてStdFSのオブジェクトを保持する。すなわち、ZFSの各ファイルについて、圧縮した内容を保存するために1つのStdFSのファイルが利用される。これにより、圧縮ファイルが作られる(図5-3:①)。ここで、StdFSの代りに、CFSのオブジェクトを下位層オブジェクトとして保持するならば、圧縮された暗号化ファイルが作られる(図5-3:②)。

オブジェクトの堆積では、堆積可能オブジェクトの順序は、固定されていない。この例では、図5-3に示すように、ZFS/CFS/StdFSの順(図5-3:②)でも、CFS/ZFS/StdFSの順(図5-3:③)でもかまわない。(ただし、圧縮と暗号の場合、順序が効率に影響を与える。先に圧縮を行う方が先に暗号化を行うよりも圧縮効率がよい。)

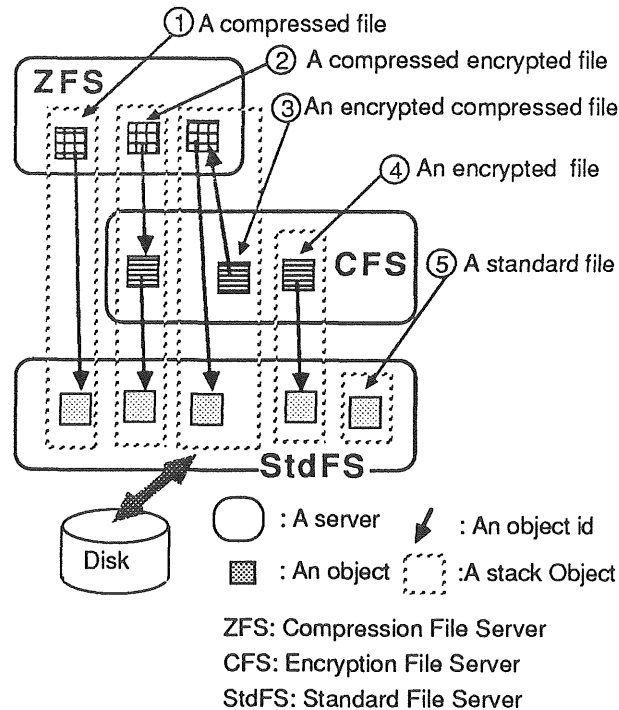


図5-3 圧縮、暗号、標準ファイル・サーバ

Figure 5-3: Compression, Encryption and Standard File Servers.

### 5. 2. 2 サーバの堆積とオブジェクトの堆積

オブジェクトの堆積では、堆積がオブジェクトを単位として作られる。一方、個々のオブジェクトではなく、サーバ全体を単位として堆積を作る方法がある。これを、サーバの堆積 (server-stacking) とよぶことにする (図5-4 (a))。サーバの堆積とオブジェクトの堆積の概念は、ほぼ同じ時期に、独立に考案された [24] [70] [85]。

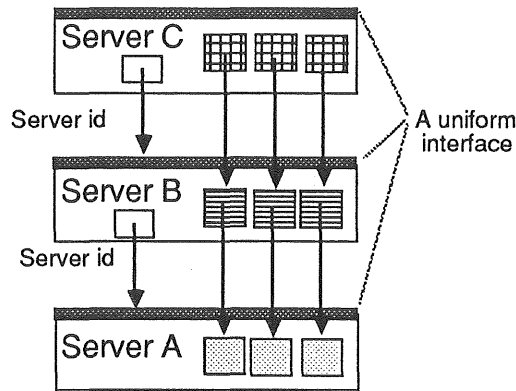


図5-4 (a) オブジェクトの概念があるサーバの堆積

Figure 5-4(a): Server-stacking with objects.

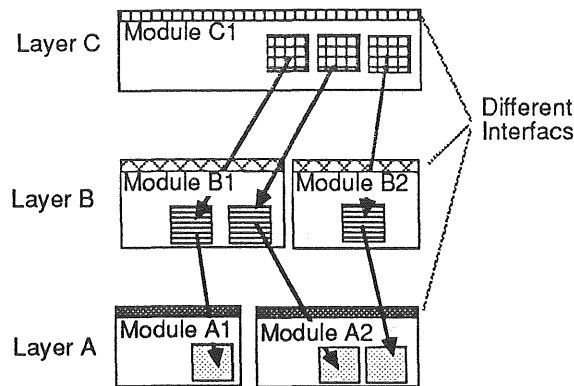


図5-4 (b) オブジェクトの概念がある層化プログラミング

Figure 5-4(b): Layered programming with objects.

図5-4 サーバの堆積と層化プログラミング

Figure 5-4: Server-stacking and layered programming.

オブジェクトの堆積では、オブジェクトを生成する時に、そのオブジェクトの下位層のオブジェクトの識別子をサーバに渡す。これに対して、サーバの堆積では、サーバを起動する時に、下位層のサーバの識別子を引数として渡す。この結果、オブジェクトの堆積と同様に、堆積可能なサーバを用いることによって、サーバの構成と順序を自由に選択することが可能となる。ここでサーバが堆積可能であるとは、自分自身

のインタフェースと、下位層のサーバのインタフェースが同一であることである。

サーバの堆積では、必ずしもサーバがオブジェクトを扱う必要はなく、単なる手続き呼出しにより機能が提供されることも考えられる。しかしながら、オブジェクトを扱うサーバの場合、サーバの堆積では、サーバ単位の自由しかなく、すべてのオブジェクトが同じ形になってしまう。図5-4(a)に、オブジェクトを扱うサーバを示す。たとえば、5.2.1項で述べたZFSとCFSをサーバの堆積により重ねた場合、すべてのファイルが暗号化された圧縮ファイルになる。暗号化だけされたファイル、あるいは、圧縮だけされたファイルを作ることができない。

オブジェクトの堆積においても、結果として、同一の形の堆積しか作られないこともある。ゆえに、サーバの堆積は、オブジェクトの堆積の保護された特殊形である。ここで保護とは、柔軟性を犠牲にして、不正な堆積が生成されないようにすることである。また、サーバ全体をオブジェクトと考えると、1つのプログラムから1つのインスタンスしか生成されないような、オブジェクトの堆積そのものであると見ることもできる。

### 5.2.3 層化プログラミングとオブジェクトの堆積

サーバの堆積は、サーバを1つの層と考えると、層化プログラミング(layered programming)と似ている(図5-4(b))。相違点は、層化プログラミングの場合、層の構成と順序が固定されているのに対して、サーバの堆積では、サーバの構成と順序が自由に選択可能になっていることである。

層化プログラミングでは、必ずしもオブジェクトの概念が必要ではない。層化プログラミングにおいて、オブジェクトの概念を導入し、1つの層について複数のサーバを実現するならば、オブジェクトの堆積と同様に、サーバを自由に組み合わせて利用することが可能になる。また、1つの層の内部の実現において、オブジェクトの堆積、または、サーバの堆積の技術を取り入れることも可能である。

### 5.2.4 UNIXのパイプラインとオブジェクトの堆積

オブジェクトの堆積の枠組みでは、UNIXにおけるパイプラインを構成する(パイプで結合された)プロセス群は、ストリームという通信プリミティブにより結合されたプロセス・オブジェクトの堆積と見ることが可能である。UNIXのパイプでは、プロセス・オブジェクトという揮発性(volatile)のオブジェクトが堆積を形成する。これに対して、ファイル・オブジェクトを積み重ねた場合、ファイルという安定(stable)なオブジェクトが堆積を形成する。たとえば、UNIXにおけるプロセス・オブジェクトの堆積では、コマンドを実行する度に新しいプロセス・オブジェクトが生成される。一方、ファイル・オブジェクトの堆積の場合、何度アクセスされても同じオブジェクトが参照される。この性質を利用してキャッシングを行うことで効率を改善することができる。

### 5. 2. 5 抽象データ型、部分型とオブジェクトの堆積

オブジェクトの堆積は、object-basedシステムだけではなく、プログラミング言語においても利用することができる。その条件は、抽象データ型 (ADT, Abstract Data Types)、および、部分型 (subtypes) を支援していることである [42]。抽象データ型とは、実現方法によって型を記述するのではなく、その型に属する値を使って利用できる操作のリスト、および、その操作の公理を示すことによって型を記述するものである。抽象データ型のデータでは、内部の状態が公開された手続きを通じてのみ操作可能である。このことは、object-basedシステムにおける遠隔手続き呼出しと対応する。

object-basedシステムにおけるオブジェクトとサーバは、プログラミング言語におけるインスタンス (オブジェクト) とクラス (抽象データ型の具体的な実現) に対応する。たとえば、ファイル抽象データ型には、read(), write() のような手続きが定義されており、すべてのファイル型のインスタンスは、この手続きを受け付けなければならない。5. 2. 1項で述べた例においては、標準ファイル・クラス、圧縮ファイル・クラス、暗号ファイル・クラスは、ファイル抽象データ型の3つ実現である。

抽象データ型には、この節の冒頭で述べた基本概念 (2) と (4) が存在しない。( (3) は、部分型と解釈される。) オブジェクトの堆積では、堆積可能オブジェクトは、インスタンス変数として、別のオブジェクトの識別子を保持し、自分自身の機能を実現するために、そのオブジェクトを呼び出さなければならない。さらに、そのオブジェクトが自分と同じ抽象データ型 (のある部分型) のインスタンス (オブジェクト) でなければならない。たとえば、堆積可能なファイル・オブジェクトは、ファイル抽象データ型の部分型のインスタンス変数を持っており、その変数が指しているオブジェクトを呼び出して、自分自身の機能を実現する。

### 5. 2. 6 C++によるオブジェクトの堆積の例

図5-5に、C++によるオブジェクトの堆積の実現の例を示す。

第6行から第10行において、Service1 は、C++の抽象クラス (abstract class) として宣言されている。C++の抽象クラスとは、抽象データ型に対応するもので、他のクラスの基底クラス (base class) となるだけで、インスタンスが生成されることがないようなクラスである。図5-5第9行の仮想関数 (virtual function) に0が代入されている。これにより、このクラス Service1 は、抽象クラスとなる。

図5-5第15行から第36行において、サービス Service1 を具体的に実現する3つのクラス ServerA, ServerB, ServerC が宣言されている。ServerA は、基底オブジェクト (bottom object) のサーバ、ServerB と ServerC は、堆積可能オブジェクトのサーバになっている。各コンストラクタは、ServerA を除き、下位層のオブジェクトの識別子を1個引数として取る。オブジェクト識別子は、C++では、オブジェクトへのポインタにより実現されている。

```

1: #include "std.h"
2:
3: //
4: //      Service
5: //
6: class Servicel                      // abstract class in C++
7: {
8: public:
9:     virtual void print( char *s ) = 0 ;
10: };
11:

```

図 5 - 5 ( a ) 抽象データ型の定義

Figure 5-5(a): The definition of the abstract data type "Servicel".

```

12: //
13: //      The interfaces of three servers.
14: //
15: class ServerA : public Servicel      // Bottom server
16: {
17: public:
18:     ServerA() {}
19:     virtual void print( char *s );
20: };
21:
22: class ServerB : public Servicel
23: {
24:     class Servicel *lower ;
25: public:
26:     ServerB( Servicel *x ) { lower = x; }
27:     virtual void print( char *s );
28: };
29:
30: class ServerC : public Servicel
31: {
32:     class Servicel *lower ;
33: public:
34:     ServerC( Servicel *x ) { lower = x; }
35:     virtual void print( char *s );
36: };
37:

```

図 5 - 5 ( b ) サーバのインタフェースの定義

Figure 5-5(a): The definitions of the interfaces of the servers.

図 5 - 5 C++言語によるオブジェクトの堆積の実現

Figure 5-5: Object-stacking in C++.

```

38: //
39: //      The implementations of the servers.
40: //
41: void
42: ServerA::print( char *s )
43: {
44:     printf("%s/A\n", s );
45: }
46:
47: void
48: ServerB::print( char *s )
49: {
50:     char buff[strlen(s)+1+2] ;
51:     sprintf( buff, "%s/B", s );
52:     lower->print( buff );
53: }
54:
55: void
56: ServerC::print( char *s )
57: {
58:     char buff[strlen(s)+1+2] ;
59:     sprintf( buff, "%s/C", s );
60:     lower->print( buff );
61: }
62:

```

図5-5(c) サーバの実現

Figure 5-5(c): The implementations of the servers.

```

63: //
64: //      The client.
65: //
66: main()
67: {
68:     Servicer *a1 = new ServerA();
69:     Servicer *b1 = new ServerB( a1 );
70:     Servicer *c1 = new ServerC( b1 );
71:     c1->print("hello[1]");
72:
73:     Servicer *a2 = new ServerA();
74:     Servicer *c2 = new ServerC( a2 );
75:     Servicer *b2 = new ServerB( c2 );
76:     b2->print("hello[2]");
77:
78:     Servicer *b3 = new ServerB( b2 );
79:     Servicer *c3 = new ServerC( b3 );
80:     c3->print("hello[3]");
81: }

```

図5-5(d) クライアント

Figure 5-5(d): The client.



```
hello[1]/C/B/A
hello[2]/B/C/A
hello[3]/C/B/B/C/A
```

図5-6 mainプログラムの実行結果

Figure 5-6: Results of the main program.

ServerB, ServerC では、各オブジェクトは、抽象データ型 `ServiceI` のインスタンスへのポインタを持っている（第24行、第32行）。これらは、下位層のオブジェクトを保持するための変数である。ServerA のオブジェクトの場合、下位層のオブジェクトを保持する変数を持っていない。

各クラスのコンストラクタは、R e S C システムにおける堆積を生成するサーバ手続き `create_with_lower` に相当する。C++では、コンストラクタ関数の名前がクラスの名前と同じという制約があるため、それぞれのクラスで別々の名前になっている。 `create_with_lower` については、5.4.2項において詳しく述べる。

第41行から第61行において、各サーバの実現が記述されている。第41行の基底オブジェクトの関数 `print` では、他のオブジェクトを呼び出すことなく、機能が実現されている。ここでは、C言語の `printf` 関数を呼び出している。第47行、および、第55行の関数 `print` では、`lower->print` という記述により、下位層のオブジェクトの機能が呼び出されている（第52行、第60行）。C++における演算子“->”によるメンバ関数の呼出しが、object-basedシステムにおける遠隔手続き呼出しに対応する。このように、堆積可能オブジェクトのメンバ関数では、下位層のオブジェクトを呼び出すことにより自分自身の機能を実現している点に特徴がある。

第66行以下、関数 `main` の部分が、上で宣言されたサーバを利用するクライアントに相当する。第68行において、サーバ `ServerA` のコンストラクタを呼び出し、基底オブジェクトを生成している。そして、生成されたオブジェクトへのポインタを `ServiceI` 型のポインタを保持する変数 `a1` に格納している。第68行では、第67行において生成されたオブジェクトへのポインタを引数として、`ServiceB` のコンストラクタが呼び出され、堆積オブジェクト `b1/a1` が生成されている。同様に、第70行では、堆積オブジェクト `c1/b1/a1` が生成されている。

第71行において、この3層のオブジェクトから構成される堆積オブジェクトの頂上オブジェクト `c1` に `print` メッセージを送っている。その結果を、図5-6の第1行に示す。このように、`c1/b1/a1` という堆積が生成されていることが分る。

図5-5の第73行から第75行では、`b2/c2/a1` という堆積オブジェクトが生成されている。このことは、図5-6の第2行の実行結果にも現れている。堆積オブジェクト `c1/b1/a1` と 堆積オブジェクト `b2/c2/a1` では、`ServerB::print()` と `ServerC::print()` の実行順序が逆になっている。このように、オブジェクトの堆積では、オブジェクトを積み重ねる順序に制限がなく、実行される手続きの順序をオブ

ジェクトごとに設定することが可能である。

図5-5第78行、第79行では、既存の堆積オブジェクト `b2/c2/a2` の上にさらに、オブジェクト `b3`, `c3` を積んでいる。その結果として、堆積オブジェクト `c3/b3/b2/c2/a2` が生成されている。第80行の `c3->print` の実行の結果を、図5-6の第3行に示す。このように、`ServerB::print` と `ServerC::print` が再帰的に呼び出されている。

この例では、C++言語の抽象データ型と部分型の機能が使われている。継承の機能が利用されているように見えるが、それは、部分型を定義するためにのみ利用されている。すなわち、部分クラスの機能は、インタフェースを一様にするためにのみ利用されており、各関数の実現を継承するためには利用されていない。

### 5.2.7 継承とオブジェクトの堆積

オブジェクトの堆積は、オブジェクト指向システムとオブジェクト指向プログラミング言語における継承 (inheritance) とは関係がない。たとえば、オブジェクト指向プログラミング言語において、5.2.1項で述べた圧縮ファイル・サーバのように、圧縮ファイル・クラスを作ることを考える。この場合、標準ファイル・クラスの部分クラスとして圧縮ファイル・クラスを開発することは、可能である。しかしながら、標準ファイル・クラスのメソッド (手続き) の実現を継承することはできない。すなわち、`read()` や `write()` といったメソッドを全て書き換える必要がある。この時、圧縮ファイル・クラスのインタフェースは、標準ファイル・クラスのインタフェースと同一のものになる。しかしながら、このことは、標準ファイル・クラスも、圧縮ファイル・クラスも、ファイル抽象データ型の1つの実現であることを示しているだけである。継承とは、インタフェースだけではなく、メソッドの実現が共有される必要がある [14]。

オブジェクト指向プログラミング言語において、継承の機能を用いることなく、5.2.1項で述べたものと同じ様なプログラムを記述することが可能である。すなわち、たまたま自分と同じ型 (の部分型) のオブジェクトの識別子をインスタンス変数として保持しており、そのオブジェクトへメッセージを送ることで、自分自身のオブジェクトの機能を実現するようなプログラムを記述すればよい。これは、オブジェクト指向プログラミング言語における抽象データ型を支援する機能を利用していることに相当する。

オブジェクト指向では、継承という機能を利用して、プログラム・コードの共有を行う。これにより、ソフトウェアの再利用可能性を高めることができる。一方、オブジェクトの堆積では、継承の機能、すなわち、メソッドのプログラムの共有は、行われない。しかしながら、インタフェースを統一することにより、汎用で独立性の高いサーバを構築することができる。これにより、プログラム・コードの共有を行うことなく、ソフトウェアの再利用可能性を高めることができる。

## 5. 2. 8 委譲とオブジェクトの堆積

オブジェクトの堆積は、オブジェクト指向プログラミング言語における委譲 (delegation) と類似点がある [48]。委譲とは、継承とは異なる仕組みにより、機能の再利用をはかるものである。委譲に基づくモデルでは、クラス概念が存在しない。新しいオブジェクトは、既存のオブジェクトをプロトタイプとして、そのプロトタイプとの違いを記述することにより作られる。オブジェクトに対してメッセージを送ると、まずそのオブジェクト自身が検索される。そして、対応するメソッドが見つかった場合、それが実行される。見つからなかった場合、そのオブジェクトのプロトタイプが呼び出される。

オブジェクトの堆積における積み重ねるという動作は、委譲におけるプロトタイプからオブジェクトを生成する動作と類似している。そして、オブジェクトの堆積における下位層のオブジェクトが、委譲におけるプロトタイプ・オブジェクトに相当する。オブジェクトの堆積と委譲との相違点を以下にまとめる。

(1) オブジェクトの堆積には、クラス (サーバ) の概念がある。オブジェクトは、クラス (サーバ) により管理される。委譲には、クラス概念が存在しない。委譲では、クラスに相当するようなオブジェクトを生成するオブジェクトを定義することも可能ではある。しかしながら、この場合においても、クラスとオブジェクトは、区別されない。

(2) 委譲には、基本概念 (4) がない。すなわち、一様なインタフェースを用いるという制約が存在しない。委譲では、一般に、オブジェクトとそのプロトタイプ・オブジェクトは、異なるインタフェースを持つ。このようなインタフェースの拡張は、継承における特殊化に相当する。

(3) オブジェクトの堆積では、未定義の手続きに対する呼出しは、下位層のオブジェクトに転送されるのではなく、サーバにより拒絶される。委譲では、自動的にそのプロトタイプに転送される。

(4) 委譲には、self という変数が存在する。その変数は、一連の手続き呼出しにおいて、最初に手続き呼出しを受け付けたオブジェクトを保持している。これに対して、オブジェクトの堆積では、そのような変数は存在しない。すなわち、上位層のオブジェクトのサーバは、一般のクライアントとして、下位層のオブジェクトのサーバを呼び出す。

## 5. 2. 9 メタ・オブジェクトとオブジェクトの堆積

オブジェクトの堆積における下位層のオブジェクトは、メタ・オブジェクトとは関係がない。メタ・オブジェクトを利用するシステムとしては、Apertoss システム (Muse システム) があげられる [109]。Apertoss は、平行オブジェクト指向計算 (concurrent object-oriented computing) とリフレクティブ計算 (reflective computing) に基づいている分散型オペレーティング・システムである。

A p e r t o s が与える単一の抽象は、オブジェクトである。各オブジェクトは、1つ、あるいは、複数のメタ・オブジェクトにより支えられている。支えているオブジェクトにより実現される環境は、メタ空間とよばれる。メタ空間は、そのオブジェクトの実行環境を提供するオブジェクトであり、インタプリタであり、そして、そのオブジェクトの実行に最適なオペレーティング・システムである。メタ・オブジェクト自身も、他のオブジェクト（メタ・メタ・オブジェクト）により、支援されている。この結果、オブジェクトは、メタ階層構造を構成する。

オブジェクトの堆積におけるサーバは、オブジェクトのメタ・オブジェクトに相当する。オブジェクトの堆積におけるサーバとA p e r t o s のオブジェクトを比較すると、オブジェクトの堆積におけるサーバは、互いに依存し合っているのに対して、A p e r t o s のオブジェクトは、因果律（the law of causality）により、依存関係が決定されている。オブジェクトの堆積では、例えば5. 2. 1項で述べたZFSとCFSは、互いに呼び出し合うことがある。これに対して、A p e r t o s では、2つのオブジェクトが互いに相手を支援し合うことはない。

オブジェクトの堆積における下位層のオブジェクトは、A p e r t o s のメタ・オブジェクトに似ている側面がある。大きな違いは、オブジェクトの堆積では、各オブジェクトの振舞がそのサーバにより決定されるのに対して、A p e r t o s では、メタ・オブジェクトにより決定される点にある。

A p e r t o s では、オブジェクトの性質を変化させる時に、別のメタ空間に移動（migrate）させることがある。たとえば、オブジェクトに永続性を与える場合、そのような機能を提供するメタ空間にオブジェクトを移動する。オブジェクトa1を、空間Aから別の空間Bに移動することができる時、AからBへ移動可能であると呼ぶことにする。この移動可能性とオブジェクトの堆積における堆積可能性は、類似点がある。すなわち、両者とも柔軟に移動／堆積するためには、一様なインタフェースが必要となる。

### 5.3 システムに対する要求

オブジェクトの堆積を実現するために必要とされるシステムの機能は、以下のよう  
にまとめられる。R e S C システムのカーネルは、これらの条件を全て満足する。

#### (1) オブジェクトの概念の存在

オブジェクトとは、データと手続きがカプセル化されたものである。オブジェクト  
のデータの参照と更新は、外部に公開された手続きによってのみ行われる。

分散型オペレーティング・システムでは、オブジェクトは、サーバにより管理され  
ている。サーバは、アドレス空間の壁により隔離されている。クライアントは、公開  
された (export された) 遠隔手続き呼出しによってのみオブジェクトの内部のデー  
タを参照したり更新したりすることができる。オブジェクト指向プログラミング言語、  
あるいは、抽象データ型を支援したプログラミング言語では、コンパイラや実行時シ  
ステムによりオブジェクトへのアクセスが公開された手続きに制限される。

R e S C は、object-based システムであり、システムの資源は、オブジェクトとし  
て提供される。オブジェクトの概念が実現ないプログラミング言語やシステムでは、  
オブジェクトの堆積を利用することができない。U N I X システムは、ファイルやプ  
ロセスをオブジェクトと考えることが可能であるので、この条件を満足する。

#### (2) オブジェクト識別子

オブジェクトは、一様なオブジェクト識別子により参照されなければならない。こ  
れは、任意のサーバのオブジェクトを下位層オブジェクトとして利用するためである。  
プログラミング言語では、オブジェクトへのポインタをオブジェクト識別子として利  
用する。この場合、型がないか、あるいは、部分型を支援している必要がある。

R e S C では、2.6 節で述べたように、システム・レベルのオブジェクトは、一  
様なオブジェクト識別子により識別される。U N I X システムにおいては、ファイル  
名をオブジェクト識別子として利用することが可能である。この場合、ファイル・オ  
ブジェクトとディレクトリ・オブジェクトを一様に扱うことができる。しかしながら、  
プロセス・オブジェクトを一様に扱うことは不可能である。また、効率の面では、文  
字列の名前よりも数値の名前の方が有利である。

#### (3) メッセージによる間接的な操作

オブジェクトの操作は、サーバに対して要求メッセージ送ることによって間接的に行われ  
る。オブジェクトをクライアントのアドレス空間にマップし、直接操作することを許  
すシステムが存在する。機能の観点から見ると、オブジェクトに対する直接操作は、  
なんら問題がない。しかしながら保護の観点からは、オブジェクトに対する不正な操

作が行われる危険性がある。プログラミング言語においては、コンパイラや実行時システムによりオブジェクトに対する不正な操作が防止される必要がある。

R e S Cカーネルは、通信プリミティブとして遠隔手続き呼出しを提供する。クライアントは、操作対象のオブジェクトの識別子を引数として、サーバに遠隔手続き呼出しを行うことで、オブジェクトを間接的に操作する。オブジェクトに対する不正な操作は、サーバにより拒絶される。

プロセス間通信機能は、非対称的でよい。すなわち、クライアントがサーバを呼び出す機能があればよく、逆に、サーバからクライアントを呼び出す機能は、不用である。

#### (4) 基底オブジェクト

オブジェクトの堆積を使っても、基底オブジェクトのサーバを作ることができないので、システムが基底オブジェクトを提供しなければならない。2. 10節で述べたように、R e S Cでは、カーネル中のサーバ群が、基底オブジェクトを提供する。プログラミング言語では、5. 2. 6項で述べたように、なんらかの方法で基底オブジェクトを実現する必要がある。5. 2. 6項の例では、C言語のライブラリ関数を用いて、基底オブジェクトを実現している。

基底オブジェクトは、単に基本的な機能を提供するだけでなく、そのインタフェースも重要である。たとえば、Vシステムのカーネルが提供しているような単純なブロック・デバイスというインタフェースでは、不十分である。Vシステムの場合、カーネル中のサーバのインタフェースが、カーネル外で構築されるファイル・サーバによる利用を前提として設計されているため、一般利用者による利用には向いていない。この場合、カーネル外の一般利用者向けのファイル・サーバのオブジェクトが、基底オブジェクトとして適している。

#### (5) 1つのサービスについて複数のサーバの存在を許す機能

R e S Cでは、オブジェクト識別子にサービス識別子ではなくサーバ識別子が含まれているので、自由にサーバを追加することができる。たとえば、ファイル・サービスやネーム・サービスのよう、システム中に既に存在しているサーバと同じサービスを提供するサーバを追加することが可能である。その追加されたサーバは、既存のシステムのサーバと全く対等に扱われる。

インターネットのサーバは、1つのサービスについて、各ホストで1つだけ走るように設計されている[101]。このサービスとサーバの1対1対応は、オブジェクトの堆積に適していない。たとえば、メール・サービスの場合、T C P / I Pのポート番号が、25番と固定されている。(U N I Xでは、この対応関係が、/etc/services というファイルに格納されている。) オブジェクトの堆積により、既存のメール・サーバの上に堆積する新たなメール・サーバを構築した場合、新たなサーバには、異なるサービス名とポート番号を割り当てなければならない。これにより、一般のク

ライアントは、両者を対等に扱うことができなくなる。

プログラミング言語において、この条件は、1つのプログラム内部で1つの抽象データ型に対して複数の実現が同時に利用可能であることに相当する。例えばC++では、5.2.6項で示したような抽象クラスを用いることにより、この条件を満足することができる。しかしながら、5.2.6項の例では、仮想関数 print の定義が各クラスのインタフェースに繰り返し現れている。このことは、一様なインタフェースを実現する妨げとなる。Eiffel言語における暫定クラス (deferred class) では、この点が改善されている [44]。

NFSが利用可能でないUNIXシステムは、この条件を満たさない [59]。すなわち、UNIXでは、カーネルが提供するファイル・オブジェクトやディレクトリ・オブジェクトの他に、同じインタフェースのファイル・オブジェクトやディレクトリ・オブジェクトを定義することはできない。NFSが利用可能なUNIXでは、利用者プロセスとしてNFSサーバを実現することが可能になっている。よって、この条件を満足する。SunOSでは、利用者プロセスによるNFSサーバを利用して、自動マウント (automount) や半透明ファイル・サービス (translucent file service) を実現している [93]。

#### (6) 下位層のオブジェクトを自由に選択する機能

ReSCのサーバは、新しいオブジェクトを作る手続きにおいて、下層のオブジェクトの識別子を引数として受け取る。基底オブジェクトの場合も、インタフェースとしては0個のオブジェクト識別子を受け取るようになっている。これにより、下位層のオブジェクトを自由に選択することができる。新しいオブジェクトを作る手続きについては、5.4.2項で詳しく述べる。

### 考察

分散核を持つ分散型オペレーティング・システムの多くは、上記(1)から(4)の条件を容易に満足する。しかしながら、条件(5)は、本質的であり、インターネットのように、この条件を満たさないために、オブジェクトの堆積が実現できないことがある。条件(6)を達成するには、それぞれの分散型オペレーティング・システムのサーバ呼出しのインタフェースを変更して、オブジェクト生成時に下位層のオブジェクトを引数として取るようにしなければならない。

## 5.4 オブジェクトの堆積による分散型オペレーティング・システムの構築

5.3節では、オブジェクトの堆積を実現する条件について述べた。この節では、分散型オペレーティング・システムを構築する際に望まれる、オブジェクトの堆積を実現する環境の性質について論じる。

#### 5. 4. 1 望まれる性質

この項では、必要条件ではないが、分散型オペレーティング・システムを構築する際に望まれる、オブジェクトの堆積を実現する環境の性質について述べる。ReSCカーネルは、ここであげる条件を全て満足する。

##### (1) 安定なオブジェクトには、安定な識別子

ファイルのような安定 (stable) なオブジェクトを扱う場合、オブジェクト識別子も安定していることが望ましい。安定とは、サーバがクラッシュして再起動した後も、変わらずに存在し続けることである。たとえば、ファイルやディレクトリは、安定なオブジェクトである。これに対して、プロセスは、揮発性 (volatile) のオブジェクトである。安定しているオブジェクトの識別子は、安定していることが望ましい。その理由は、サーバがクラッシュして再起動した場合、クラッシュの前とオブジェクト識別子が変わらない方が好ましいからである。これに対して、プロセスのように揮発性のオブジェクトの場合、揮発性のオブジェクト識別子で十分である。

2. 6 節で述べた ReSC のオブジェクト識別子の中で、システム・レベルで定義されているサイト識別子とサーバ識別子は、安定である。オブジェクト識別子の安定性は、各サーバにおけるオブジェクト番号の管理方法に依存している。

##### (2) コネクションレスの通信プリミティブ

通信プリミティブとしては、コネクションレスのものが望ましい。その理由は、上記 (1) と同様に、サーバのクラッシュと再起動に耐えられるからである。コネクション付きの場合、サーバの再起動により、コネクションを維持することができない。ReSC でも、コネクションレスの通信プリミティブである遠隔手続き呼出しを採用している。

##### (3) 位置依存のオブジェクトと位置独立のオブジェクト

オブジェクトの堆積を使って、分散処理を行うサーバを構築する場合、位置独立のオブジェクトと位置依存のオブジェクトの両方を扱う必要がある。たとえば、ファイルの複製では、1 つのファイルの複製を複数のサイトに配置する。そして、あるサイトに障害が発生しても、別のサイト上にある複製 (replica) により、そのファイルへのアクセスが可能になる。このことは、位置に依存したオブジェクトを生成する必要があることを示している。逆に、自分自身のオブジェクトは、位置独立でなければならない。それは、せっかく複製を作っても、自分自身が位置依存では、そのサイトが落ちた場合、複製が働かないからである。

2. 6 節で述べたように、ReSC では、位置依存のオブジェクトと位置独立のオブジェクトの両方が識別可能な識別子を利用している。分散処理を行うサーバについては、5. 6 節において述べる。



#### (4) 位置独立のプロセス間通信機能

オブジェクトの堆積に基づき、分散型オペレーティング・システムを構築する場合、オブジェクトの位置に独立したプロセス間通信機能が不可欠である（単にオブジェクトの堆積を実現するためには、この機能は不用である。）。2.7節で述べたように、ReSCでは、位置依存なオブジェクト識別子により指し示されているオブジェクトに対しても、位置に依存せず遠隔手続き呼出しを行うことができる。

#### (5) 軽量プロセス

軽量プロセスは、サーバ・プロセスのように、複数のクライアントを同時に扱う必要があるプログラムを記述する際に、有効な方法であることが知られている[98]。それは、独立したクライアントからの要求を並列に処理することができるからである。単一プロセッサにおいても、軽量プロセスを利用することにより、複数のクライアントとの通信処理が並列に実行される。

3章で述べたように、ReSCは、マイクロプロセスと仮想プロセッサという概念に基づく軽量プロセス機能を提供する。軽量プロセスを用いたサーバの構成については、5.9.1項で述べる。

#### (6) 要求メッセージの転送 (forward)

Vシステムのように、遠隔手続き呼出しを転送 (forward) する機能があると便利である。その理由は、オブジェクトの堆積では、何もしないで下位層オブジェクトへ要求メッセージを伝えることが多いからである。この場合、転送機能を利用することで、応答メッセージが直接最初のクライアントに戻されるので、メッセージの受渡しのオーバーヘッドが軽減される。さらに、サーバを軽量プロセスで実現する場合、同時に活動する軽量プロセスの数を少なくすることができる。その理由は、下位層オブジェクトからの応答を待つ必要がなくなるからである。2.7節で述べたように、ReSCカーネルは、要求メッセージの転送機能を持っている。

#### (7) 分離されたネーム・サーバ

ネーム・サーバ (ディレクトリ・サーバ) は、オブジェクトの管理と独立している方が望ましい。その理由は、オブジェクトの堆積では、サーバの種類が非常に多くなるからである。仮にサーバごとに名前空間を構成したとすると、オブジェクトのサーバの種類ごとに名前空間が構成され、利用者インタフェースの質が低下する。

ReSCのネーム・サーバも、オブジェクトのサーバと分離している。

### 5.4.2 堆積を生成する手続き

オブジェクトの堆積では、クライアントが堆積の構造を意識して利用することも、意識しないで利用することも可能にすることが望まれる。ReSCでは、各サーバに

次の3つのサーバ手続きの定義を義務付けることにより、これを実現している。

(1) `oid_t create_with_lower( oid_t lowers[], opt_t options[] )`

これは、オブジェクトの堆積を直接的に実現する手続きである。この手続きは、与えられたオブジェクト識別子で指定されたオブジェクトを下位層オブジェクトとするような新たなオブジェクトを生成し、生成したオブジェクトの識別子を返す。

(2) `oid_t copy( oid_t old )`

これは、与えられた引数 `old` で示されたオブジェクトに対して、いわゆる「深いコピー (deep copy)」を行う事で、新しいオブジェクトを生成する手続きである。下位層オブジェクトがある場合、それらは、再帰的にコピーされる。

(3) `oid_t create_like( oid_t old )`

これは、手続き (2) と同じく、引数 `oid` で示されたオブジェクトと同じ形のオブジェクトを生成する手続きである。(2) との違いは、オブジェクトの内容がコピーされない点にある。結果として、空のオブジェクトが生成される。

手続き (1) を使うことにより、クライアントは、堆積の構造を意識してオブジェクトを生成することができる。手続き (2)、または、手続き (3) を使うことにより、クライアントは、そのオブジェクトの堆積の構造を気にすることなく、同じ構造を持つ新しいオブジェクトを生成することができる。手続き (2) は、コピー・コマンドで使われる。手続き (3) は、エディタなどで、編集しているファイルと同じ堆積の構造を持つファイルを作るときに使われる。

#### 5. 4. 3 ひな形オブジェクト・ベース

R e S C では、頻繁に利用されるオブジェクトを予め生成し、あるディレクトリに登録しておく。このようなディレクトリを、ひな形オブジェクト・ベース (template object bases) とよぶ。たとえば、次の様な名前のファイルが登録されている。

(1) `std-file`

(2) `compressed-std-file`

(3) `encrypted-compressed-std-file`

上記 (1) は、5. 2. 1 節で述べた標準ファイル・サーバ `S t d F S` の空のオブジェクトである。すなわち、`S t d F S` の空のファイルが "std-file" という文字列の名前でそのディレクトリに登録されている。(2) は、`Z F S` のオブジェクトで、その下位層オブジェクトが `S t d F S` の空のオブジェクトである。(3) は、`Z F S` の

オブジェクトで、下位層オブジェクトとしてCFSのオブジェクトを持っている。ここで、そのCFSのオブジェクトは、下位層オブジェクトとしてStdFSのオブジェクトを持っている。

ここに登録されているオブジェクトは、5.4.2節で述べたサーバ手続き `create_with_lower()` を利用し、システム管理者がオブジェクトを1個1個積み上げることで作られたものである。一方、一般の利用者がこれらのオブジェクトを使う場合、必要なオブジェクトをコピー・コマンドでコピーして利用する。そのコピー・コマンドでは、サーバ手続き `copy()` が利用されている。

#### 5.4.4 書込み時コピー

MaChでは、書込み時コピー (copy-on-write) という技術を使って、UNIXのプロセス生成システム・コール `fork` を効率的に実現している [66]。

`fork` システム・コールでは、親プロセスのメモリ・イメージを子プロセスへコピーする。書込み時コピーでは、実際にコピーを行うことを、書込みが起こるまで遅らせることで、書込みが起こらないページが親子間で共有される。そのため、親プロセスと子プロセスの両方の読み書き可能領域に、書込み保護を掛けておく。どちらかのプロセスがその領域に対して書込みを行うと、ページ・フォルトが発生する。カーネルは、ページ・フォルトの処理の中で、新しいページを割り当て、そのページの内容をコピーし、以後、親プロセスと子プロセスが別々のページを参照するようにアドレス変換表を再設定して、実行を再開する。書込みが行われない限り、ページは、親プロセスと子プロセスの間で共有される。この結果、`fork` システム・コールの効率が改善される。

ReSCでは、5.4.2節で述べた `copy()` 手続きの実現において、この手法を利用することができる。この場合、プロセス・オブジェクトだけでなく、ファイルやディレクトリにおいても利用することができる。たとえば、あるファイルがコピーされ、元のファイルに対しても新たに生成されたファイルに対しても書込みが行われなかった場合、それらのファイルの実体が同一のディスク・ブロックを共有するような実現を行うことが可能となる。

#### 5.5 種々の堆積可能オブジェクトとサーバ

この節では、オブジェクトの堆積に基づき高度な機能を提供するオブジェクトとそのサーバについて述べる。ここで述べるサーバのオブジェクトは、堆積可能である。すなわち、それらのサーバのオブジェクトを積み重ねることで、それらのサーバの機能を統合して利用することができる。ReSCシステムでは、このような高度な機能は、主に逐次応用プログラムにより利用される。

ここで述べるサーバのいくつかは、R e S C以外の分散型オペレーティング・システムでも利用可能であろう。ただし、5. 4節で述べたような機能がない場合、一部の機能を利用することができない。

## 分類の観点

堆積可能なオブジェクトとサーバは、次のような観点から分類される。

### (1) 機能

クライアントからメッセージを受け取り、対応する下位層のオブジェクトにメッセージを送る時に、メッセージの内容になんらかの操作をほどこすものをフィルタ・オブジェクト (filter object) とよぶ。これに対して、ほとんど変換を行わないようなオブジェクトを間接オブジェクト (indirection object) とよぶ。そして、それぞれのオブジェクトのサーバをフィルタ・サーバ (filter server)、および、間接サーバ (indirection server) とよぶ。

### (2) 対象となるオブジェクトの型

下位層のオブジェクトとして、例えば、ファイルという特定の型に対して動作するオブジェクトをファイル・オブジェクトとよぶ。これに対して、下位層に任意の型のオブジェクトを取ることができるオブジェクトを総称オブジェクト (generic object) とよぶ。そして、それぞれのサーバをファイル・サーバ、総称サーバ (generic server) とよぶ。

### (3) 各オブジェクトごとの下位層のオブジェクトの数

1つのオブジェクトが下位層のオブジェクトとして複数のオブジェクトを利用するものをグループ・オブジェクト (group object) とよぶ。そして、それを管理するサーバをグループ・サーバ (group server) とよぶ。

### (4) 分散処理を行うかどうか

下位層のオブジェクトが異なるサイト (ネットワーク上に分散したプロセッサ) にあることを利用して分散処理を行うオブジェクトを分散型オブジェクト (distributed object) とよぶ。そして、それを管理するサーバを分散型サーバ (distributed server) とよぶ。

この節では、上記(1)から(3)について、それぞれ代表的なサーバの機能と内部動作を説明する。(4)については、5. 6節で述べる。

### 5.5.1 フィルタ・サーバ

フィルタ・サーバとは、受け付けた要求を下位層オブジェクトへ中継する際に、さまざまな変換処理を行うサーバである。

#### 5.5.1.1 フィルタ・ファイル・サーバ

UNIXでは、多くのコマンドがパイプにより結合することで組み合わせて利用することが可能になっている。このようなコマンドは、フィルタ型コマンド (filter) と呼ばれる。フィルタ型コマンドの例としては、tr, sort, sed があげられる。

UNIXにおける各フィルタ型コマンドから、堆積可能サーバを構築することができる。たとえば、次のUNIXのコマンド行：

```
% tr A-Z a-z <from >to
```

に対応して、ファイル中の大文字を小文字に変換するファイル・サーバを構築することができる (図5-7 (a))。このように、ファイルを対象とし、フィルタ型コマンドと同等の機能を提供するサーバを、フィルタ・ファイル・サーバ (filter file server) と呼ぶことにする。5.2.1項で述べたZFSとCFSも、フィルタ・ファイル・サーバに分類される。

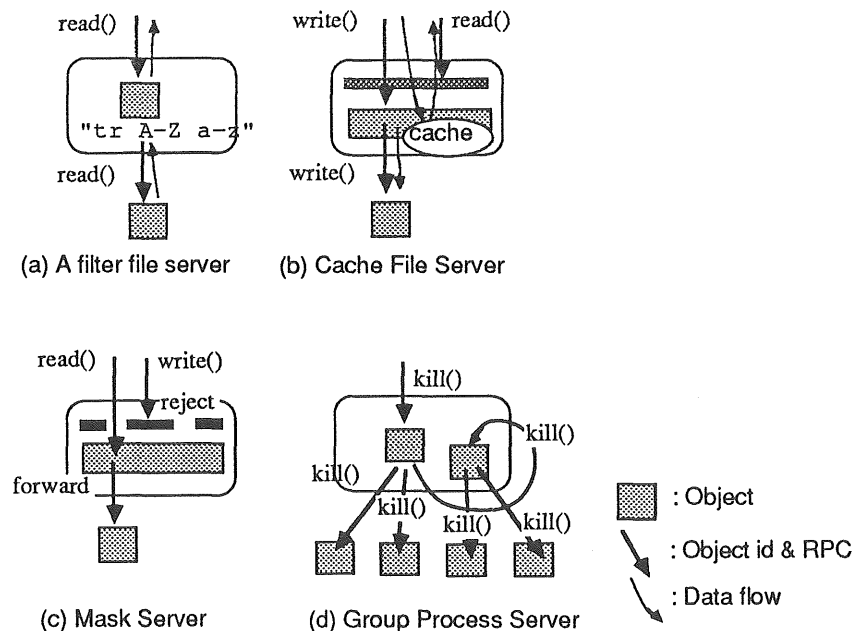


図5-7 さまざまな堆積可能サーバ

Figure 5-7: Various stackable servers.

一般のフィルタ型コマンドには、逆関数が存在しないので、それから構築されるファイル・サーバは、読み専用（または、書き専用）になる。5. 2. 1項で述べた圧縮ファイル・サーバと暗号ファイル・サーバでは、compress, encrypt に対して、uncompress, decrypt という逆関数が存在するので、読み書き可能なファイル・サーバを構築することが可能である。

#### 5. 5. 1. 2 フィルタ・ディレクトリ・サーバ

ファイルと同様に、ディレクトリに対してもフィルタを考えることができる。たとえば、あるディレクトリについて、“\*. [ch]” というパターンにマッチするファイルだけを含むようなディレクトリを作ることが考えられる。

#### 5. 5. 1. 3 フィルタ・ウインドウ・サーバ

ウインドウ・システムでは、次の様なウインドウ・サーバを構築することが考えられる。

- (1) スケールを変える。
- (2) 回転する。
- (3) 色を変える。

このようなサーバを組み合わせることにより、このような機能を持たないウインドウに対して、上記の機能を付加することができる。

### 5. 5. 2 間接サーバ

間接サーバとは、受け付けたメッセージの内容を変更することなくそのまま下位層オブジェクトへ転送するようなサーバである。この時、単純に転送するだけでなく、2次的な操作を行う事で、様々な機能を提供することができる。

#### 5. 5. 2. 1 キャッシュ・ファイル・サーバ

キャッシュ・ファイル・サーバとは、オブジェクトに対する要求を受け付けると、対応する下位層オブジェクトに問い合わせをせず自分自身で答えるようなサーバである。これを、図5-7 (b) に示す。クライアントから書き込み要求を受け付けると、それをキャッシュに書き込み、下位層オブジェクトに転送する。読み込み要求を受け付けると、まずキャッシュを調べ、ヒットした場合はそれを返す。ミスした場合、下位層オブジェクトに読み込み要求を送り、その内容を自分自身のキャッシュに書き込みを行った後、クライアントに返す。

このキャッシュ・ファイル・サーバの特徴は、どのようなファイル・オブジェクトについてもキャッシングを行うことができる点にある。堆積可能ではないファイル・サーバでは、そのサーバが管理するファイルについてしかキャッシングを行うことができない。

#### 5. 5. 2. 2 キャッシュ・ディレクトリ・サーバ

キャッシュ・ファイル・サーバと同様に、ディレクトリ・オブジェクトについても、キャッシュ・ディレクトリ・サーバを構築することができる。

#### 5. 5. 2. 3 マスク・サーバ

マスク・サーバとは、オブジェクトに送られてきた要求の中で特定の要求だけを下位層オブジェクトへ転送し、そのほかの要求をさえぎるようなサーバである。たとえば、図5-7(c)では、読み込み要求をそのまま転送し、書き込み要求をさえぎることで、書き込み禁止ファイルが作られている。

マスク・サーバは、総称サーバ (generic server) である。すなわち、このサーバは、任意の型のオブジェクトを扱うことができる。たとえば、ディレクトリ・オブジェクトに対しても、ファイル・オブジェクトに対しても、マスク・サーバを利用することができる。

#### 5. 5. 3 グループ・サーバ

上記のサーバでは、各オブジェクトについて下位層オブジェクトを1個取るものであった。しかしながら、オブジェクトの堆積において、下位層オブジェクトの数が1個に制限されているわけではない。すなわち、複数の下位層オブジェクトを取ることにより、オブジェクトのグループを作ることができる [75]。

オブジェクトのグループ化の基本的な考え方は、次の通りである。

- (1) 1つのオブジェクトに対して、複数の下位層オブジェクトを保持する。
- (2) 要求を受付けると、それを下位層オブジェクトへ多重送信 (multiplex) する。
- (3) 下位層オブジェクトからの応答を多重受信 (demultiplex) し、クライアントへ応答する。

#### 5. 5. 3. 1 ディレクトリの融合

ディレクトリの融合とは、複数のディレクトリの内容を全て含むような1つのディレクトリを作成することである。堆積可能なグループ・サーバの1つとして、ディレクトリの融合を行うサーバを構築することができる。この時、融合されるディレクトリは、単純なディレクトリでも既に融合されたディレクトリでもよい。すなわち、階層的にディレクトリの融合を行うことができる。

5. 9. 2項において、キャッシングの技術を用いたディレクトリの融合を行うサーバの実現について述べる。

#### 5. 5. 3. 2 ファイルの結合

ファイルの結合とは、複数のファイルを順に結合したような1つのファイルを作成することである。ファイルの結合を行うサーバは、5. 5. 1. 1で述べたフィルタ・フ

ファイル・サーバにおいて、UNIXの cat コマンド (concatenate) に対応したものととして構築される。

### 5. 5. 3. 3 プロセス・グループ

プロセス・グループとは、そのグループに属する複数のプロセスを同時に操作するための仕組みである。オブジェクトの堆積では、堆積可能なプロセス・グループを実現することが可能である。図5-7 (d) では、そのプロセス・グループに属するプロセスに対して、プロセスを殺す操作 (kill) が行われている。

Vシステム [12] は、カーネルにおいてプロセス・グループの機能を提供している。Vシステムでは、プロセス・グループを、マルチキャスト通信を実現するために利用している。Vシステムのプロセス・グループは、堆積可能ではない。すなわち、他のプロセスの実現やプロセス・グループ自身をメンバとするグループを構築することができない。

### 5. 5. 3. 4 ウィンドウの飾り付け

ウィンドウ・システムでは、ウィンドウ・マネージャにおいてタイトル・バー、スクロール・バーなどを付加する、ウィンドウの飾り付け (decoration) を実現している。これは、オブジェクトの堆積の枠組みでは、ウィンドウの本体、タイトル・バー、スクロール・バーを下位層オブジェクトとしてもつウィンドウ・オブジェクトを生成することに相当する。こうして生成されたウィンドウ・オブジェクトは、堆積可能であり、他の堆積可能なウィンドウ・オブジェクトと組み合わせて利用することが可能である。

## 5. 6 分散型堆積可能オブジェクトとサーバ

この節では、分散型堆積可能オブジェクトとそれを管理するサーバについて述べる。分散型堆積可能オブジェクト (distributed stackable object) とは、5. 5 節の冒頭で述べたように、下位層のオブジェクトが異なるサイトにあることを利用して分散処理を行うようなオブジェクトである。そして、それを管理するサーバを分散型堆積可能サーバ (distributed stackable server) とよぶことにする。

### 5. 6. 1 複製ファイル・サーバ

複製ファイル・サーバ (Replication File Server) とは、ファイルの複製を行い、フォールト・トレラントなファイル・サービスを提供するサーバである。複製ファイル・サーバは、1つのオブジェクトについて、複数のサイト上の複数の下位層オブジェクトを利用する。これを、図5-8に示す。このサーバは、書込み要求を受け付けると、全ての下位層オブジェクトへ送る。読み込み要求を受け付けると、いくつかの下



位層のオブジェクトに要求を送り、そのうち少なくとも1個のオブジェクトから応答が得られれば十分である。このことにより、ある下位層オブジェクトが存在するサイトが落ちている場合も、他の下位層のオブジェクトの働きにより、そのファイルにアクセスすることが可能となる。

この複製ファイル・サーバの動きから、複製とは、次の様な意味であると考えることができる。

- (1) 複数のサイトに、同一の内容を持つオブジェクトを生成する。
- (2) クライアントに対して、それらがあたかも1つのオブジェクトであるかのごとく見せる。
- (3) 具体的には、受け付けたメッセージを下位層オブジェクトへ多重送信 (multiplex) し、下位層オブジェクトからの応答を多重受信 (demultiplex) することで実現することができる。

オブジェクトの複製では、下位層オブジェクトの位置が重要である。たとえば、あるファイル・オブジェクトの全ての下位層オブジェクトが同一サイトにある場合、複製が無効になる。したがって、下位層のオブジェクトを指し示すために、2.6節で述べた位置依存のオブジェクト識別子を用いる。これとは対照的に、複製サーバ自身のオブジェクトを指し示すためには、位置独立のオブジェクト識別子を用いる。

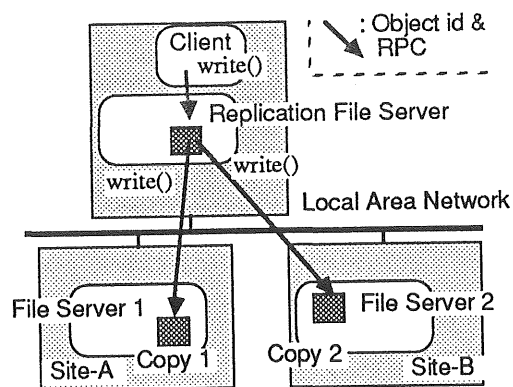


図5-8 複製ファイル・サーバ

Figure 5-8: Replication File Server

### 5.6.2 間接オブジェクトによるオブジェクトの移動

オブジェクトの移動 (object migration) は、凍結 (freeze)、解凍 (unfreeze)、下位層の変更 (change\_lower) という3つの手続きを持つ間接オブジェクトを利用して実現することができる (図5-9)。

このオブジェクトは、通常、受け付けた要求を自ら解釈することなく、下位層のオブジェクトへ転送 (forward) する。たとえば、読み込み要求や書き込み要求を下位層へ

転送する。しかしながら、凍結 (freeze)、解凍 (unfreeze)、下位層の変更 (change\_lower) という3つの手続きだけは、自ら解釈する。オブジェクトの移動は、次のような手順で実現される。

- (1) オブジェクトの移動の管理者 (initiator) は、移動するオブジェクト (間接オブジェクト) に凍結要求 (freeze) を送る。それ以降、そのオブジェクトへの通常のメッセージは、下位層のオブジェクトへ転送される代わりに、間接サーバの中の待ち行列に保持される。
- (2) 管理者は、下位層のオブジェクトの内容を新しいオブジェクトにコピーする。
- (3) 管理者は、新しい下位層のオブジェクトの識別子を引数として、下位層の変更要求 (change\_lower) を間接オブジェクトに送る。間接オブジェクトのサーバは、そのオブジェクトの下位層のオブジェクトを新しいものに置き換える (図5-9)。
- (4) 管理者は、その間接オブジェクトに解凍要求 (unfreeze) を送る。この時、そのオブジェクトの待ち行列に要求が保持されていた場合、間接オブジェクトのサーバは、それらの要求を新しい下位層のオブジェクトへ転送する。

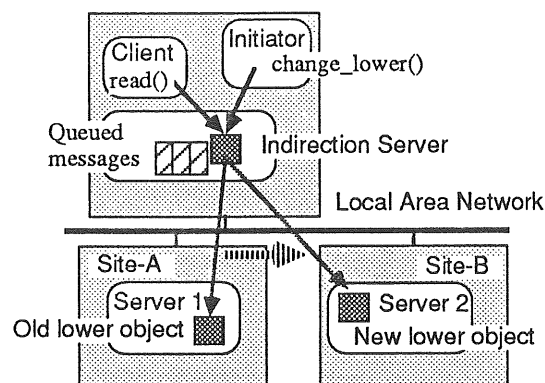


図5-9 間接サーバによるオブジェクトの移動

Figure 5-9: Object migration by an indirection server.

この間接サーバの動きより、オブジェクトの移動とは、次のような意味であると考えられることができる。

- (1) 元のオブジェクトが存在するサイトとは別のサイト上に新しいオブジェクトを作る。
- (2) クライアントに対して、2つのオブジェクトがあたかも同一のものであるように思わせる。
- (3) 具体的には、オブジェクトの移動は、オブジェクトのアクセスを全て1回間接で行う事で実現される。

ここで説明した間接オブジェクトとそのサーバは、5.5節で述べた総称オブジェクト、および、総称サーバである。この間接オブジェクトは、ファイル・オブジェクト

トだけでなく、ディレクトリ・オブジェクトの移動にも利用することができる。

ここで説明した間接オブジェクトとそのサーバは、分散型堆積可能オブジェクト、および、分散型堆積可能サーバである。この間接サーバの場合、同じ時刻には、1つのサイト上のオブジェクトしか利用しない。しかしながら、異なる時刻においては、複数のサイト上のオブジェクトを利用することがある。

### 5. 6. 3 分散型堆積可能オブジェクト間の競合

分散型堆積可能オブジェクトは、しばしば互いに競合する。たとえば、5. 6. 2項で述べた間接オブジェクトを下位層のオブジェクトとして、5. 6. 1で述べた複製ファイル・オブジェクトを生成することを考える。この場合、間接オブジェクトの働きで下位層のオブジェクトが全て1つのサイトに移動した場合、複製が意味をなさなくなる。

分散型堆積可能オブジェクト間の競合は、従来の分散型オペレーティング・システムにおいてマッピング・コントローラを含む並列／分散応用プログラムを走らせた時に生じる、システムと並列／分散応用プログラムとの競合と同じものである。今後の研究の発展の1つとして、この競合を解消し調和をはかる方法の研究があげられる。

## 5. 7 多重視点

5. 5. 3項で述べたオブジェクトのグループ化では、複数の下位層オブジェクトを管理するオブジェクトを利用した。逆に、1つの下位層オブジェクトが複数の上位オブジェクトから参照されることも考えられる。これは、オブジェクトの多重視点 (multiple views) を実現していることに相当する。たとえば、図5-10では、1つの下位層のファイル・オブジェクトについて、単純な間接オブジェクト、書込みをマスクするオブジェクト、大文字を小文字に変換するフィルタ・ファイル・オブジェクトを利用することにより、次の3つの視点を得られる。

(視点1) 読み書き可能なファイル

(視点2) 読み専用ファイル

(視点3) 小文字だけからなるファイル

このように、オブジェクトの堆積を使うことにより、多重視点を直接的に表現することが可能となる。

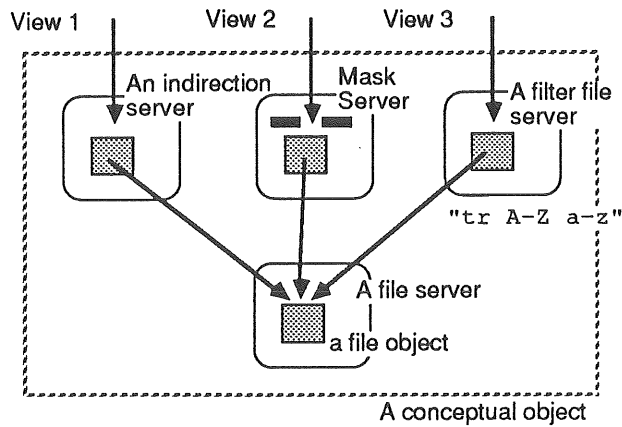


図5-10 3つの上位層のオブジェクトによる3つの視点の実現

Figure 5-10: An object with three views by three upper objects.

### 5.7.1 堆積の上下を参照することによる多重視点

全ての堆積オブジェクトは、途中の層をアクセスすることにより多重視点を実現していると考えられる。図5-11に、3つのオブジェクトから構成されている堆積オブジェクトの3つの視点を示す。

むやみに堆積の途中のアクセスを許した場合、問題が生じることがある。たとえば、マスク・サーバを用いて、あるファイルを読み込み専用にしたとしても、その下位層オブジェクトへのアクセスを許せば、読み込み専用にした意味が失われる。このような問題を未然に防ぐためには、下位層オブジェクトのオブジェクト識別子を秘密にする必要がある。このことは、Amoebaシステム [54] [55] におけるケーパビリティと同様に、暗号の技術を用いることで実現可能である。

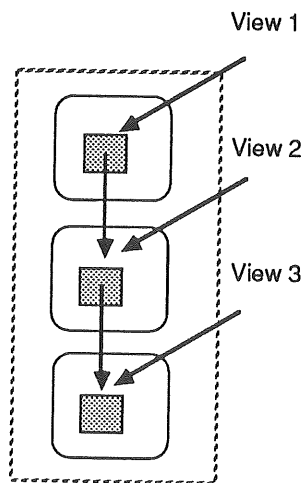


図5-11 堆積オブジェクトによる3つの視点の実現

Figure 5-11: A stack object with three views.

### 5. 7. 2 多重視点とキャッシュ

多重視点を実現する場合、問題となることの1つは、キャッシュの一貫性である [107]。たとえば、図5-10の例において、個々の視点ごとに5.5.2.1で述べたキャッシュ・オブジェクトを積み重ねることを考える。この時、読み書き可能ファイルの視点に対して書込みが行われた場合、他の視点のキャッシュの一貫性が保たれなくなる可能性がある。

この問題の解決方法の1つは、同期型キャッシュ・ファイル・サーバ (coherent cache file server) を構築することである。このサーバでは、図5-12に示すように、すべての視点のキャッシュが関連付けられて管理される。そして、ある視点に書込みが行われた場合、他の視点のキャッシュが無効化される。

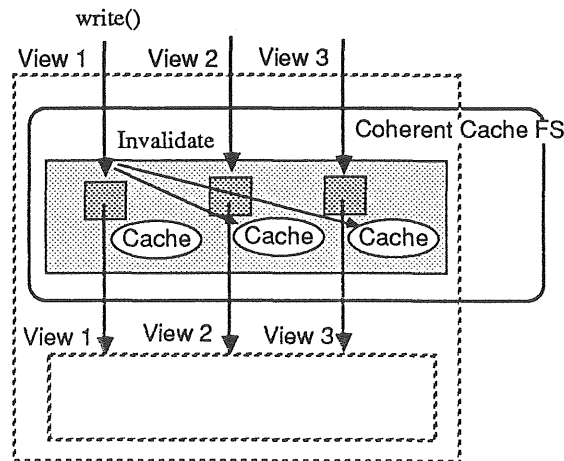


図5-12 同期型キャッシュ・サーバ

Figure 5-12: Coherent Cache File Server.

## 5. 8 遠隔手続き呼出しスタブ生成器とインタフェース記述言語

分散型オペレーティング・システムにおいて、オブジェクトの堆積を利用する場合、普通、システムにより提供される遠隔手続き呼出しが利用される。多くの分散型オペレーティング・システムでは、遠隔手続き呼出しのインタフェースを記述するための言語と、その言語から遠隔手続き呼出しのスタブを自動的に生成するプログラムを備えている。このプログラムは、スタブ生成器 (stub generator)、あるいは、スタブ・コンパイラと呼ばれている。このようなインタフェース記述言語やスタブ生成器の例としては、Am o e b a の A I L (Amoeba Interface Language) [54]、M a c h システムの M i G (Mach Interface Generator) [105] [58]、S u n R P C の r p c g e n (RPC protocol compiler) [69] の記述言語があげられる。

従来のスタブ生成器のインタフェース記述言語は、サーバを中心にインタフェースを記述するものであった。このことは、オブジェクトの堆積を利用する場合、一様なオブジェクトのインタフェースの記述を困難にする。そこで、オブジェクトの堆積に適したインタフェース記述言語の設計を行った。この節では、設計したインタフェース記述言語の構文と意味について述べる。そして、従来のインタフェース記述言語として

S u n R P C のインタフェース記述言語をとりあげ、比較を行う。

### 5. 8. 1 インタフェース記述言語の構文と意味

本インタフェース記述言語において扱えるデータ型は、C言語、および、S u n R P C のインタフェース記述言語 [69] を拡張したものになっている。拡張した型構成子は、proc, sproc, set of, server である。付録B-1に、本インタフェース記述言語の詳細を示す。

本記述言語において、インタフェースは、次のような文の並びにより記述される。

<名前> : <型> = <値> ;

この構文は、f u n [9] に準じている。f u n は、1階型付きλ計算に基づく、多相 (polymorphism) とオブジェクト指向言語をモデル化するために設計された言語である。ただし、f u n において関数を意味する "fun" と記述するところを、本インタフェース記述言語では、手続きを意味する "proc", あるいは、"sproc" と記述する。また、f u n では、型と値を明確に区別しているのに対して、本インタフェース記述言語では、型を型型の値として扱う。

図5-13に示した簡単なファイル・サーバの記述を例に、この言語の特徴について述べる。

```

1: buff_t: type = array of char ;
2:
3: kill:   proc( void ) void           = 14 ;
4: create: sproc( lower:array of oid_t ) new:oid_t = 1 ;
5: copy:   sproc( orig:oid_t ) new:oid_t      = 2 ;
6:
7: any_server: set of sproc = { copy, create };
8: any_object: set of  proc = { kill };
9:
10: read:  proc( where,count:int ) buff_t = 52 ;
11: write: proc( where,count:int;buff:buff_t ) void = 53 ;
12:
13: file: set of proc = any_object + { read,write };
14: read_only_file: set of proc = file - { write };
15: file_server: type = server( any_server + file );
16:
17: StdFS: file_server = 3 ;
18: ZFS  : file_server = 5 ;
19: CFS  : file_server = 6 ;
20: ReadOnlyFS : server( any_server+read_only_file ) = 7 ;

```

図 5 - 1 3 簡単なファイル・サーバのインタフェース記述

Figure 5-13: The interface description of simple file servers.

型構成子 `proc` は、オブジェクト手続きを、型構成子 `sproc` は、サーバ手続きを定義する。オブジェクト手続き／サーバ手続きとは、5. 2 節の冒頭で説明したように、それぞれオブジェクト／サーバに対する遠隔手続き呼出しのエントリである。それぞれは、オブジェクト指向言語におけるインスタンス・メソッド／クラス・メソッドに相当するものである。たとえば、第 10 行では、手続き `read` が 2 つの整数型の値を入力パラメタとし、バッファ型 (`buff_t`) の値を出力パラメタとするオブジェクト手続きとして定義されている。

オブジェクト手続きは、隠れたパラメタとしてオブジェクト識別子 (`oid_t`) を含むようなサーバ手続きである。たとえば、第 10 行に定義されているオブジェクト手続き `read` は、次の様なサーバ手続きと等価である。

```
read: sproc( this:oid_t; where,count:int ) buff_t = 52 ;
```

オブジェクト手続き型とサーバ手続き型の変数の値 (たとえば、`read` では 52) は、遠隔手続き呼出しの手続き番号である。この値は、システム全体で一意でなければならない。スタブ生成器は、この一意性を確認する。

この言語の特徴は、型構成子 `set of` を用いて手続きの集合としてインタフェースを記述する点にある。たとえば、第 7 行では、`any_server` という変数にサーバ手続き `create` とサーバ手続き `copy` を要素とする集合が代入されている。これは、任意のサーバがサーバ手続き `create` とサーバ手続き `copy` をインタフェースに含んでいることを意味している。

この言語では、手続きの集合を扱う演算として、集合の和と集合の差がある。たと

例えば、第13行では、任意のオブジェクトに対して、読み込み、および、書き込みを行う手続きを追加して、普通のファイル・オブジェクトのインタフェースが定義されている。次の行では、読み込み専用のファイル・オブジェクトのインタフェースが、普通のファイル・オブジェクトから書き込みを行う手続きを削除することによって定義されている。

型構成子 `server` は、手続きの集合を引数としてサーバ（サーバ識別子型）を定義する。第15行では、ファイル・サーバが、任意のサーバとしての手続きとファイルの読み書きを行う手続きを受け付けることができるサーバとして定義されている。（オブジェクト手続きは、オブジェクト識別子を隠れた引数として持つような特殊なサーバ手続きであることから、サーバ手続きの集合とオブジェクト手続きの集合の和を求めることができる。） `ReSC` では、この型の変数の値は、サーバ識別子として利用されるので、全体で唯一でなければならない。スタブ生成器は、その唯一性を確認する。

このスタブ生成器は、以下の点においてオブジェクトの堆積に適している。

- (1) サーバではなく、オブジェクトを中心にインタフェースを定義することが可能となっている。このことは、オブジェクトを中心に考えるオブジェクトの堆積に合致している。
- (2) オブジェクトの堆積において重要となる一様なインタフェースを持つオブジェクトを容易に定義することができる。
- (3) オブジェクト手続きが、特殊なサーバ手続きであるという `object-based` システムの性質を直接的に反映した記述が可能となっている。

図5-13の第4行において定義されているサーバ手続き `create` が、オブジェクトを積み重ねる操作を行うものである。これは、5.4.2項で述べた `create_with_lower` を簡略化したものである。この図の `create` は、引数としてオブジェクト識別子の配列を受け取り、それらを下位層のオブジェクトとして新しいオブジェクトを生成し、その識別子を返すサーバ手続きとして定義されている。

### 5.8.2 SunRPCのインタフェース記述言語との比較

図5-13に示したファイル・サーバの、`SunRPC`におけるインタフェースの記述を、図5-14に示す。図5-13に示した本記述言語では、第17行と第18行において、`StdFS` と `ZFS` が同一のインタフェースを持っていることが明示的に示されている。これに対して、図5-14に示した `SunRPC` の記述言語では、`StdFS` と `ZFS` の記述が第18行と第28行に散在している。さらに、それぞれのサーバにおいて定義されている手続きのインタフェースが同一であることを保証する手段が存在しない。



```

1: struct oid_t {
2:     site_t      site_id;
3:     server_t    server_id;
4:     object_number_t object_number;
5: };
6:
7: typedef opaque buff_t<>;
8: struct create_arg_t
9: {
10:     oid_t  lower<>;
11: };
12: struct read_arg_t {
13:     oid_t  this ;
14:     int    where;
15:     int    count ;
16: };
17:
18: program StdFS_PROGRAM {
19:     version StdFS_VERSION {
20:         void_t  KILL(oid_t)           = 14 ;
21:         oid_t   CREATE(create_arg_t)  = 1 ;
22:         oid_t   COPY(oid_t)           = 2 ;
23:         buff_t  READ(read_arg_t)      = 55 ;
24:         void_t  WRITE(write_arg_t)    = 56 ;
25:     } = 1;
26: } = 3 ;
27:
28: program ZFS_PROGRAM {
29:     version ZFS_VERSION {
30:         void_t  KILL(oid_t)           = 14 ;
31:         oid_t   CREATE(create_arg_t)  = 1 ;
32:         oid_t   COPY(oid_t)           = 2 ;
33:         buff_t  READ(read_arg_t)      = 55 ;
34:         void_t  WRITE(write_arg_t)    = 56 ;
35:     } = 1;
36: } = 5 ;

```

図5-14 SunRPCにおける簡単なファイル・サーバのインタフェース記述

Figure 5-14: The interface description of the simple file servers in SunRPC

本記述言語のその他の特徴としては、引数と結果の数が可変であることがあげられる。rpcgenでは、入力パラメタと出力パラメタの数がそれぞれ1個に制限されている。たとえば、図5-13の第10行における手続きreadの入力パラメタに対応するものは、図5-14の第12行で定義されている構造体read\_arg\_tにまとめられている。図5-14の第13行のthisは、5. 8. 1項で述べたオブジェクト手続きの隠されたパラメタである。

## 5.9 堆積可能サーバの実現

この節では、典型的な堆積可能サーバの構造と実現方法を示し、その実現を支援するために構築したライブラリの機能について述べる。ここで示す堆積可能サーバの実現において特徴的なことは、キャッシングの技術を利用している点、および、3章で述べた軽量プロセスを利用している点にある。

堆積可能サーバとライブラリは、現在、R e S Cカーネルではなく、既存のシステム上に構築したR e S Cカーネルの機能をエミュレートする環境において利用可能となっている。この節では、この環境についても述べる。

### 5.9.1 典型的な堆積可能サーバの構成要素

典型的な堆積可能サーバは、次のようなモジュールから構成されている(図5-15)。

- ・ 遠隔手続き呼出し受け付け軽量プロセス
- ・ 遠隔手続き呼出し入口手続き
- ・ メモリ中のオブジェクトに対する手続き
- ・ 局所オブジェクト管理部

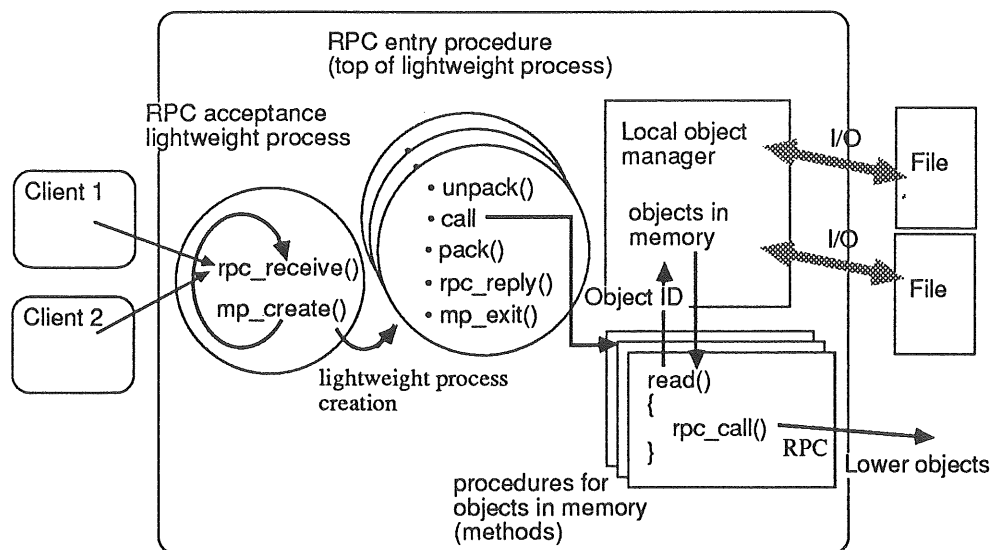


図5-15 典型的な堆積可能サーバの構造

Figure 5-15: The structure of a typical stackable server.

#### ■ 遠隔手続き呼出し受け付け軽量プロセス

遠隔手続き呼出し受け付け軽量プロセス (RPC acceptance lightweight process) とは、クライアントからの遠隔手続き呼出しの受け付けを行う軽量プロセスである。この

軽量プロセスは、通常、クライアントからの遠隔手続き呼出しを受け付けるカーネル・コールを発行している状態にある。遠隔手続き呼出しを受け付けると、手続き番号を確認し、以下で述べる遠隔手続き呼出し入口手続きをトップとする軽量プロセスを生成する。その手続きが存在しなければ、クライアントへエラーを返す。

#### ■遠隔手続き呼出し入口手続き

遠隔手続き呼出し入口手続き (RPC entry procedure) とは、遠隔手続き呼出しのセッションに対応して最初に呼出される手続きである。これは、軽量プロセスのトップとして、遠隔手続き呼出し受け付け軽量プロセスから呼び出される。

遠隔手続き呼出し入口手続きでは、次のような処理を行う。

- (1) クライアントから1個にまとめられて送られてきた遠隔手続き呼出しの引数を解く (unpack, unmarshall)。
- (2) 対応するメモリ中のオブジェクトに対する手続きを呼び出す。
- (3) (2) 手続きの結果を1個にまとめ (pack, marshall)、遠隔手続き呼出しの結果としてクライアントに返す (rpc\_reply())。
- (4) 軽量プロセスを終了する。

#### ■メモリ中のオブジェクトに対する手続き

メモリ中のオブジェクトに対する手続き (procedures for objects in memory) とは、オブジェクト指向プログラミング言語、あるいは、抽象データ型を支援しているプログラミング言語における、手続き (メソッド) に対応するものである。通常のプログラミング言語における手続きと比較して、特徴的な処理としては、次のようなものがある。

- (1) アクセス権のチェック

サーバとして、遠隔手続き呼出しを行ったクライアントがそのオブジェクトに対してその手続きを実行する権利があるかどうかを調べる。

- (2) 相互排除

複数のクライアントが1つのオブジェクトに対して同時に遠隔手続き呼出しを発行することがある。並列に実行することができない遠隔手続き呼出しの場合、ロック等を用いて相互排除を行う。相互排除は、軽量プロセス間の同期プリミティブにより実現される。

#### ■局所オブジェクト管理部

局所オブジェクト管理部 (local object manager) とは、メモリ中のオブジェクトを管理するモジュールである。このモジュールは、次の様な機能を提供する。

- (1) 新たにオブジェクトを生成する場合、オブジェクト識別子の割当てを行う。
- (2) 与えられたオブジェクト識別子に対応するオブジェクトのデータをファイルから読み込み、メモリ中のオブジェクトに変換する。ファイルに格納されるオ

プロジェクトのデータを、ファイル構造体 (file structure) とよぶ。

(3) メモリ中のオブジェクトをファイルへ構造体として保存する。

(4) メモリ中のオブジェクトのキャッシュを管理する。

このモジュールは、UNIXのアイノード (inode) とディスク・アイノードの管理とよく似ている。このモジュールにおけるオブジェクト識別子、メモリ中のオブジェクト、ファイル構造体が、UNIXにおけるアイノード番号、アイノード、ディスク・アイノード (disk inode) に対応する。UNIXとの違いを、以下にまとめる。

- ・UNIXでは、アイノードもアイノードから参照されるデータ (ファイルのデータ) も、ディスク・ブロックに格納される。これに対して、本モジュールでは、ディスク・ブロックではなく、ファイルに格納される。
- ・UNIXでは、アイノード番号は、カーネル内にあり、利用者から直接利用されることはない。これに対して本モジュールが扱うオブジェクト識別子は、利用者により操作される。これにより、利用者によるオブジェクト識別子の改変 (forging) や偽造 (tampering) に対処する必要がある。

図5-16に、典型的な堆積可能サーバにおけるオブジェクト識別子の構造と、2つのファイルを用いたオブジェクトの格納の様子を示す。メモリ中のオブジェクトのデータは、固定長部分 (fixed length part) と可変長部分 (variable length part) に分けられ、別々のファイルに保存される。オブジェクト識別子のオブジェクト番号の構造は、個々のサーバにより決定される。一般には、次の3つ組から構成される。

< シリアル番号, ノード番号, チェック >

シリアル番号は、オブジェクト識別子の一意性を実現するためのものである。ノード番号は、UNIXのアイノード番号と同様に、固定長部分を格納するファイル中の構造体の位置を示すインデックスとして利用される。最後のチェックは、オブジェクト識別子の改変と偽造を防ぐためのものである。この部分を利用して、Amoebaシステム [54] [55] におけるケーパビリティで用いられている技術と同様の技術により、利用者によるオブジェクト識別子の改変と偽造に対処する。

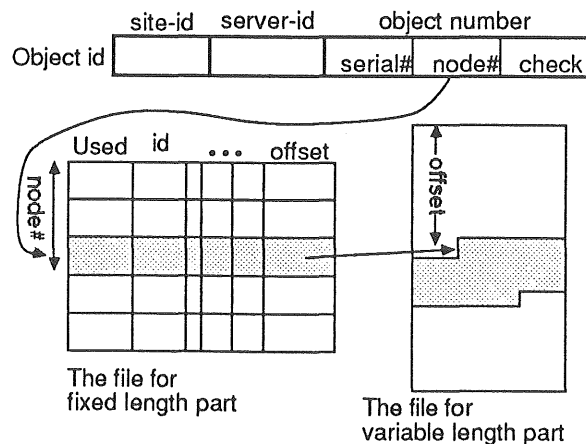


図5-16 オブジェクト識別子とファイル構造体

Figure 5-16: An object id and its file structure.

以下に、与えられたオブジェクト識別子をメモリ中のオブジェクトに変換する手順を示す。

- (1) オブジェクト番号の中のチェックを利用して、オブジェクト識別子が破損していないかどうかを調べる。破損していた場合は、エラーを返す。
- (2) メモリ中のキャッシュを検索する。ヒットした場合は、それを返す。
- (3) オブジェクト番号の中のノード番号から、オブジェクトのデータの固定長部分を読み込む。固定長部分には、オブジェクト識別子が含まれている。これと現在変換中のオブジェクト識別子を比較して、等しいことを確認する。等しくない場合は、エラーを返す。
- (4) 必要ならば、オブジェクトのデータの可変長部分を読み込む。
- (5) ファイルに格納されないオブジェクトのデータを初期化する。
- (6) オブジェクトをキャッシュに登録する。

### 5. 9. 2 キャッシングを利用した融合ディレクトリ・サーバの実現

この項では、5. 5. 3. 1において述べたディレクトリの融合を行うサーバを例に、キャッシングを利用した堆積可能サーバの構築について述べる。ここでは、特徴的な手続きである `create_with_lower`, `dir_readdir`, `dir_lookup` を取上げる。

`create_with_lower`: この手続きは、5. 4. 2項で説明したように、下位層のオブジェクトとなるオブジェクトの識別子を引数として、新しいオブジェクトを生成するものである。この手続きが呼び出されると、サーバは、与えられたオブジェクトがディレクトリかどうかを確認する。局所オブジェクト管理部を呼び出して、オブジェクト識別子を割り当て、メモリ中にオブジェクトを生成する。最後に、割り当てられたオブジェクトの識別子をクライアントに返す。

`dir_readdir`: この手続きは、ディレクトリの一覧を返すものである。まず、そのディレクトリに対するキャッシュがメモリ中に存在するかどうか調べる。存在しない場合、下位層のオブジェクトに `dir_readdir` メッセージを送り、それらの結果を融合し、結果をキャッシュとしてメモリ中に保存する。キャッシュが存在する場合、下位層のオブジェクトの最終更新時刻を調べ、キャッシュの整合性を確認する。キャッシュが古くなっていた場合、それを再構成する。最後に、キャッシュの内容をクライアントに返す。

`dir_lookup`: この手続きは、文字列の名前を受け取り、対応するオブジェクト識別子を返すものである。まず、`dir_readdir` と同様に、整合性の取れたキャッシュを作る。次に、与えられた文字列をキーとしてキャッシュを検索する。最後に、見つかったエントリのオブジェクトの識別子をクライアントに返す。

### 5. 9. 3 外部のオブジェクトへのキャッシュの保存

前項（5. 9. 2 項）で述べたディレクトリの融合の場合、扱う対象がディレクトリであるため、全てのキャッシュをサーバのメモリ中に置くことが可能である。しかしながら、大きなデータを必要とする場合、キャッシュをサーバのメモリ中に格納することは困難になる。この場合、キャッシュをサーバの外部のオブジェクトに格納する。

図5-17は、5. 2. 1 項で述べた圧縮ファイル・サーバ（ZFS）の、キャッシュを用いた実現を示している。この実現では、解凍したデータをキャッシュとしてサーバの外部のファイル・オブジェクトに格納している。ZFSのオブジェクトは、クライアントから読み込み、あるいは、書込み要求を受け付けると、下位層のオブジェクトに読み込み要求を送り、圧縮されたデータを読み込む。次に、そのデータを解凍し、外部のファイル・オブジェクト（一時ファイル（temporary file））に書き込む。そして、受け付けた読み込み要求や書込み要求を一時ファイルへ転送（forward）する。クライアントからのアクセスが途絶えた場合、ZFSサーバは、一時ファイルの内容を読み込み、圧縮し、対応する下位層のオブジェクトへ書き戻す。

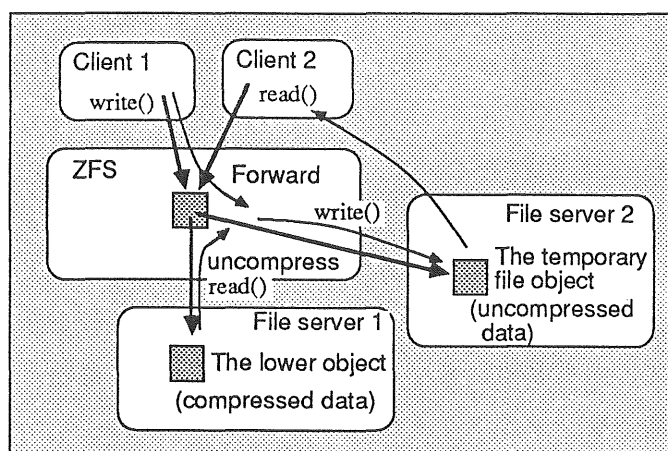


図5-17 外部のオブジェクトへキャッシュを保存するZFSの実現

Figure 5-17: The implementation of ZFS by using external objects as caches.

### 5. 9. 4 オブジェクトの堆積を実現する環境

オブジェクトの堆積を実現する環境を、SunRPCとNFSを利用して実現した。2章で述べたRescの遠隔手続き呼出しを、SunRPCを用いてエミュレートした。基底オブジェクトを提供するファイル・サーバとして、NFSサーバを直接呼び出すものを実現した。

この環境では、Rescのオブジェクト識別子は、次の様に実現されている。

- ・サイト識別子：  
インターネット・ホスト・アドレス（4バイト）を利用した。
- ・サーバ識別子：  
各サーバについて、4バイトの小さな整数を割り当てた。
- ・オブジェクト番号：  
32バイト長の不定型（opaque）（バイト列）を利用した。これは、NFSにおけるオブジェクト識別子である file handle と同じ長さである [59]。これにより、オブジェクト番号に NFS の file handle を格納することで、NFSサーバを直接呼び出すような基底ファイル・サーバの実現が容易になっている。

SunRPCのライブラリ関数では、静的変数を利用している。したがって、3.6.4項で述べたように、そのままでは、マイクロプロセス・ライブラリと共存できないものがある。そこで、SunRPCのソース・コード確認して、共存不可能なものを削除した。

SunRPCでは、スタブ生成器（rpcgen）により、サーバの主プログラムが自動生成される。しかしながら、この主プログラムは、軽量プロセスを利用するようにはなっていない。これに代わるものとして、5.9.1項で述べた、遠隔手続き呼出し受け軽量プロセスと遠隔手続き呼出し入口手続きをライブラリ関数として用意した。サーバの開発者は、従来のSunRPCのサーバ側の手続きに相当する、メモリ中のオブジェクトに対する手続きを開発する。そして、その手続きの表を、そのライブラリ関数に渡すことで、サーバを記述する。

局所オブジェクト管理部は、各サーバにより個別に記述されるものである。しかしながら、オブジェクトのデータをファイルに格納する部分には、共通な操作が多い。そこで、それを行うものを、ライブラリ関数として用意することにした。構造体のファイルへの格納には、SunRPCにおける標準データ表現形式であるXDR（eXternal Data Representation）を利用した。

## 5.10 関連した研究

5.2節では、オブジェクトの堆積のモデルとして特徴を、サーバの堆積、層化プログラミング、UNIXのパイプ、抽象データ型、継承、委譲、メタ・オブジェクトとの比較しながら述べた。この節では、その他の関連した研究との比較を行う。

### 5. 10. 1 P l a n 9

P l a n 9 は、A T & T ベル研究所で開発されているオペレーティング・システムである。P l a n 9 では、ディレクトリの融合を行う機能をカーネルにおいて提供している。本システムでは、カーネルの外の堆積可能サーバにより、ディレクトリの融合を実現している。P l a n 9 のディレクトリの融合の機能は、堆積可能ではない。すなわち、カーネルが扱う 1 種類のディレクトリについてのみ融合を支援している。堆積可能サーバでは、様々な種類のディレクトリを融合することが可能になっている。

P l a n 9 では、ウインドウ・システムのインタフェースとして、ファイル・サーバと同一のインタフェースを利用している。これにより、ファイル进行操作するプリミティブと同一のプリミティブを用いて、ウインドウ进行操作することが可能になっている。

P l a n 9 では、ウインドウ・システムの開発において、サーバの堆積と似た手法を利用している。新たにウインドウ・システムを開発する場合、まず、既存のウインドウ・サーバ上で動作するウインドウ・サーバを開発する。この瞬間には、2 つのバージョンのウインドウ・サーバが堆積している。このような環境において、ウインドウ・システムのクライアントを開発する。その後、新しいウインドウ・サーバを単独で起動する。ただし、サーバの堆積を利用しているのは、ウインドウ・システムの開発の局面だけであり、5. 5. 1. 3 で述べたフィルタ・ウインドウ・サーバのようなサーバは、利用されていない。

### 5. 10. 2 interposition

文献 [ 3 3 ] [ 3 4 ] では、interposition という概念が提案されている。interposition とは、仮想計算機、あるいは、5. 2. 9 項で説明したメタ空間に近い概念である。interposition では、プロセスが発行したシステム・コールやオブジェクトに対するアクセスを、捕捉 (trap) し、プロセスのサーバへ伝える。サーバでは、捕捉したシステム・コールやオブジェクトの呼出しに変換を加えたり、あるいは、チェックを行い、再び、別のシステムやサーバを呼び出す。この時、このシステムやサーバの呼び出しが、さらに別のサーバにより捕捉され、変換が加えられる可能性がある。

5. 2. 9 項において述べたように、メタ空間の階層とオブジェクトの層は、互いに直交した概念である。同様に、interposition とオブジェクトの堆積も直交した概念である。たとえば、オブジェクトの堆積では、複数のサーバを同時に利用することができるが、interposition では、複数の仮想計算機を同時に利用することはできない。一方、interposition では、全てのシステム・コールやオブジェクトへの参照を捕捉することが可能であるが、オブジェクトの堆積では、不可能である。

### 5. 10. 3 多重継承

オブジェクト指向プログラミング言語では、複数のクラスが提供する機能を同時に



利用するために、多重継承 (multiple inheritance) が用いられる。多重継承では、メソッドを組み合わせる方法 (メソッド結合) が重要な研究課題となっている。Flaver では、before、after、primary 等の指定を行い、多重継承におけるメソッドの実行順序を決定している [104]。この方法は、非常に複雑でわかりにくいものであった。新Flaver では、利用者によるメソッド結合の定義が可能になっている。

オブジェクトの堆積における、下位層のオブジェクトの同名の手続きを呼び出すことが、(多重)継承におけるスーパークラスの同名のメソッドを呼び出すことに対応する。オブジェクトの堆積では、個々の手続きの実現において、手続き呼出しの順序だけでなく、呼出しそのものを行うかどうかを自由に制御することができる。また、同名の手続きを呼び出さず、他の名前の手続きを呼び出すことで、その手続きを実現することもある。例えば、5.9.2項で述べた融合ディレクトリ・サーバの dir\_lookup 手続きの実現では、下位層の dir\_lookup 手続きを呼び出すのではなく、dir\_readdir を呼び出している。このように、オブジェクトの堆積では、一般的なメソッド結合のパターンを定義することはできないが、個々の手続きにおいて、自由に手続きを組み合わせることが可能になっている。

実時間 Mach のスケジューラの実現では、オブジェクトの堆積と類似の技術が利用されている [56] [102]。実時間システムにおけるスケジューラを構築する場合、様々なスケジューリングの方針 (policy) を組み合わせる必要がある。実時間 Mach では、方針を実現するものが、一様な外部仕様を持つオブジェクトとして定義されている。利用者は、それらのオブジェクトをいくつかのリンクを用いて結合することで、それらのオブジェクトが提供する機能を組み合わせて、実時間スケジューラを構築する。

この方法は、文献 [56] では、多重継承の柔軟な実現方法として提案されている。しかしながら、その論文で提案されている方法は、5.2.7項で行った継承の定義に従えば、継承ではなく、オブジェクトの堆積に分類される。すなわち、継承では、インタフェースだけでなく、メソッド (手続き) の実体 (プログラム) が共有されることが必要である。文献 [56] で提案されている方法では、インタフェースが一様になるが、メソッドは、共有されない。

#### 5.10.4 手続きの堆積

オブジェクトの堆積では、オブジェクト識別子、あるいは、オブジェクトへのポインタを用いて、複数のオブジェクトを結合し、それらのオブジェクトが持つ機能を組み合わせて利用する。オブジェクトの堆積では、オブジェクト、および、オブジェクト識別子という概念が必須である。これに対して、オブジェクトやオブジェクト識別子を一切用いることなく、手続きを積み重ねて、複数の手続きの機能を統合して利用する方法がある。これを手続きの堆積 (procedure-stacking) とよぶことにする。

手続きの堆積を支援しているプログラミング言語として、Forth がある [20] [31]。Forthでは、手続き（Forthでは、ワードと呼ばれる）の中で同名の手続きを呼び出した場合、再帰呼出しではなく、直前に定義された同名の手続きの呼出しとなる。次の例では、(2)において W1 を呼び出しているが、それは、自分自身の呼出し（再帰呼出し）ではなく、(1)において定義されている W1 の呼出しになる。

```
: W1 <処理 1> ; \ (1)
: W1 <処理 2> W1 <処理 3> ; \ (2)
```

ここで、コロン ":" に続く文字列が手続きの名前である。それ以降、セミコロン ";" までが手続きの処理を表す。バックスラッシュ "\ " 以降は、コメントである。上記の定義は、次の定義と等価である。

```
: W1 <処理 2> <処理 1> <処理 3> ; \ (3)
```

手続きの定義中で同名の手続きを呼び出した場合に再帰呼出しではなく既に定義されている別の手続きを呼び出すという意味付けは、オブジェクト指向プログラミング言語における差分プログラミングと同等の機能を提供する。たとえば、次の例では、C言語の構文を用いて、マルチプロセッサ用にロック機能を付加した1文字出力手続きを定義している。

```
putchar( c )
{
    lock( &putchar_lock );
    putchar( c ); /* call the non-lock version of putchar */
    unlock( &putchar_lock );
}
```

この例では、putchar 手続き（C言語では、関数）にロックの機能を付加している。

オブジェクトの堆積と手続きの堆積を比較すると、オブジェクトの堆積では、オブジェクト識別子を使って手続きを呼び出す順序を制御している点に特徴がある。これにより、同じサーバが管理するオブジェクトでも、オブジェクトごとに組み合わせられる手続きの種類と順序が異なる。

## 5. 1 1 まとめ

オブジェクトの堆積とは、object-based システムにおいて、複数のオブジェクト

が持つ機能を統合して利用するためのモデル化の手法である。オブジェクトの堆積は、非常に単純なモデルであるが、非常に強力であり、応用範囲も広い。この章では、オブジェクトの背景となる技術と基本概念について述べた。重要な基本概念は、積み重ねること、および、一様なインタフェースであった。次に、オブジェクトの堆積による分散型オペレーティング・システムの構築について述べた。オブジェクトの堆積により、キャッシングやオブジェクトの移動といった高度な機能が提供される。これらのサーバの実現の特徴は、キャッシングの技術が利用されている点にある。関連した研究として、サーバの堆積、層化プログラミング、継承、委譲、メタ・オブジェクト等を取りあげ、比較を行った。このようにして、オブジェクトの堆積の特徴と有効性を明らかにした。

## 第6章 結論

本研究では、並列／分散応用プログラム、および、逐次応用プログラムを対象としたオペレーティング・システムの設計と実現を行った。第1章では、研究の背景、目標、および、関連した研究について述べた。

並列／分散応用プログラムの特徴は、内部に複数の処理単位と資源割当てモジュールを含んでいる点にある。これに対して、逐次応用プログラムは、処理単位を1つだけ含み、資源割り当てモジュールを含んでいない。従来のオペレーティング・システムは、主に集中型システムで開発された逐次応用プログラムを想定してシステムの機能が設計されていた。このため、並列／分散プログラムという、内部に資源割当てモジュールを含むプログラムへの対応が不十分であった。具体的には、オペレーティング・システム内部の資源割当てモジュールの方針と並列／分散応用プログラムの資源割当てモジュールの方針が競合することがあった。これをシステムと並列／分散応用プログラム間の競合とよんだ。また、並列／分散応用プログラムにとっては、オペレーティング・システムの資源割当てモジュールが不用であるが、逐次プログラムにとっては、必要となる。これを逐次応用プログラムと並列／分散応用プログラム間の競合とよんだ。

本研究では、逐次応用プログラムに加えて、並列／分散応用プログラムという、内部に資源割当てモジュールを含むプログラムを対象とした分散型オペレーティング・システムの設計と実現を行った。このシステムをR e S Cと名付けた。R e S Cが対象とするハードウェア環境は、共有メモリ型マルチプロセッサと単一プロセッサが高速LANにより結合されたものである。

R e S Cが対象とする応用プログラムは、並列／分散応用プログラム、および、逐次応用プログラムである。

本研究の目的を、上記の2種類の競合を解消し、次の2種類の調和を実現することに設定した。

(1) オペレーティング・システムと並列／分散応用プログラム間の調和：

オペレーティング・システムは、並列／分散応用プログラムと競合することなく、そのような応用プログラムに対して、高速化や高信頼性を実現するためにマッピングやスケジューリングを制御する機能を提供する。

(2) 逐次応用プログラムと並列／分散応用プログラム間の調和：

これらの応用プログラムが、必要に応じてシステムの機能を利用することができるようにする。システムを集中型システムとも分散型システムとも見ることが可能にする。

この目標を設定したことは、従来の分散型オペレーティング・システムの定義と研究

目標を設定し直すことを意味する。従来の分散型オペレーティング・システムは、分散透明性が実現され、仮想的な集中型システムというモデルを提供するオペレーティング・システムとして定義されていた。そして、分散透明性を実現することが研究の目標であった。本研究では、並列／分散応用プログラムを考慮して、従来の定義と目標であった透明性の実現を、調和の実現へ改めた。

上記の調和を実現するために、本研究では、従来の分散型オペレーティング・システムを、カーネル上で走る並列／分散応用プログラムとして構成する方法を用いた。これは、次の3つの意味を持つ。

(1) 選択可能化：

各応用プログラムがシステムの機能を選択的に利用することを可能にする。

(2) 隔離：

競合する主体間の情報交換を制限し、独立に動作させる。

(3) オペレーティング・システムの応用プログラム化：

システムの構成要素と応用プログラムが対等な立場で計算機資源を利用する。オペレーティング・システムを並列／分散応用プログラムと同一の技術を用いて構築する。

本論文の第2章から第5章では、様々な局面で現れる競合に対する調和を図る方法を示した。

第2章では、本システムの全体構成とカーネルの機能について述べた。本システムは、カーネル、外部サーバ、並列シェル、ライブラリから構成される。カーネル化されたカーネルと比較して、本システムの特徴は、効率を重視する並列／分散応用プログラムのために、カーネル中に完全なサーバを置いている点にあった。また、システム・レベルのオブジェクト識別子は、位置独立のオブジェクトと位置依存のオブジェクトの両者を識別することが可能になっている。このオブジェクト識別子と、位置独立の遠隔手続き呼出しの発行により、分散型ハードウェア環境を、分散型システムとも集中型システムとも見ることが可能になることを示した。

カーネルは、オペレーティング・システムの中で唯一、応用プログラムによる選択ができない部分である。本研究では、選択可能化を進めるにあたり、並列／分散応用プログラムがしばしば自ら制御することを好む機能をカーネル外の構成要素で提供することにした。同時に、公平な資源割当て、応用プログラム間の通信、および、保護を実現する機能をカーネルに残した。これにより、オペレーティング・システムにおけるカーネルの論理的な位置付けを明らかにした。

第3章では、軽量プロセスの実現方式について述べた。軽量プロセスは、多重プログラミング・システムにおける、各応用プログラム内部の並列処理の単位としてのプロセスである。本研究では、マイクロプロセスと仮想プロセッサという概念を用いて軽量プロセスを実現する方式を提案した。これは、従来のカーネル制御方式と比較し

て、より効率的な実現が可能になっている。また、コルーチン方式と比較して、並列処理が可能になっているという利点がある。

本軽量プロセス実現方式の特徴は、応用固有の軽量プロセス間の同期・通信プリミティブの開発、および、スケジューラの開発が容易である点にある。これは、軽量プロセス機能が層構造を持つライブラリにより実現されていることによる。カーネル・レベルの仮想プロセッサにより、利用者レベルの軽量プロセスを実行する方式は、本研究の他にも提案された。これらの方式と比較して、本研究の特徴として、まず、層構造のライブラリにより利用者レベルの軽量プロセスが実現されていることがあげられる。次に、本仮想プロセッサの機能の特徴は、各仮想プロセッサごとに固有のメモリ領域が存在することである。これにより、利用者レベルの軽量プロセスの実現が容易になり、かつ、並列処理の効率を改善することが可能となる。

本研究では、仮想プロセッサ機能を提供するカーネルの構成法を提案した。この構成法の特徴は、カーネル内の軽量プロセスを、固有のメモリ領域を持つ実プロセッサにより実行する点にある。この方式は、R e S C システムの仮想プロセッサだけではなく、他の外部仕様を持つ仮想プロセッサの実現においても利用可能である。この方式では、カーネルの一部を利用者プロセスとして開発することが可能となる。これにより、カーネルの動作の確認とデバッグが容易になる。

実現した軽量プロセスと他の軽量プロセス、および、既存のオペレーティング・システムの重量プロセスとの基本的な性能を比較する実験を行った。また、データベースの並列処理システムを例として、応用固有の同期・通信プリミティブの開発、および、スケジューラの開発を行い、開発したシステムの性能を調べる実験を行った。これらの実現と実験を通じて、提案した軽量プロセス実現方式の有効性を示した。

軽量プロセスの実現についても、選択可能化が適用されている。それは、層構造を持つライブラリに端的に現れている。一方、応用プログラムによる選択が不可能なカーネルには、公平なCPU資源を実現するプロセス・スケジューラを置いた。システムのスケジューラと複数の応用プログラムのスケジューラは、小数のカーネル・コールによってのみ結合されている。こうしてスケジューラを互いに隔離することにより、それらの間の競合を防止している。また、仮想プロセッサを提供するカーネルの構成法には、オペレーティング・システムの応用プログラム化が適用されている。すなわち、カーネル自身をカーネル上で動作する並列プログラムと同一の技術で開発している。結果として、軽量プロセスに関して、選択可能化、隔離、および、オペレーティング・システムの応用プログラム化が有効であることが、確認された。

第4章では、マッピング・コントローラについて述べた。マッピング・コントローラとは、並列／分散応用プログラムの中に在り、ネットワーク上のサイト（プロセッサ）に対して、プロセスやデータの配置の方針を決定するモジュールである。本研究では、カーネルには全くマッピング・コントローラを含めず、各応用プログラムごとにマッピング・コントローラを開発する方針を提案した。これにより、並列／分散プ

プログラムとシステムの間でのマッピングに関する競合が解消される。システムは、応用固有のマッピング・コントローラの開発を支援するために、資源の位置や資源の利用状況の情報を積極的に提供する。

マッピング・コントローラを持たない逐次応用プログラムに対しては、システムは、並列シェルを提供する。並列シェルは、互いに協調して動作するコマンド群を並列に実行するためのコマンド・インタプリタである。並列シェルは、マッピング・コントローラを含まない応用プログラムに代わって、マッピングの最適化を行う。

ネットワークで結合されたワークステーション上に並列シェルのプロトタイプを実現した。そして、実現した並列シェルを利用して、並列処理の実験を行った。その結果、プロセス間の結合関係を利用することで、並列処理の効率が改善されることが確認された。

マッピングにおいて、カーネルからマッピング・コントローラの機能を並列シェルに移動した点に、選択可能化が適用されている。一方、それらのマッピング・コントローラは、互いに隔離されており、プロセスやファイルを生成することで資源の利用状況を変化させたり、カーネルに資源の利用状況を問い合わせたりする以外に、情報交換を行わない。また、システムの構成要素である並列シェルが、他の並列／分散応用プログラムと対等な立場で動作する点に、オペレーティング・システムの応用プログラム化が適用されている。結果として、マッピングにおいて、選択可能化、隔離、および、オペレーティング・システムの応用プログラム化が有効であることが確認された。

第5章では、オブジェクトの堆積について述べた。オブジェクトの堆積とは、object-basedシステムにおいて、複数のオブジェクトの機能を組み合わせて利用するためのモデル化の手法である。ReSCシステムでは、主に逐次プログラムに対して高度な機能を提供するために、外部サーバを統合する場面においてオブジェクトの堆積が利用される。第5章では、オブジェクトの堆積の基本概念、それを利用する条件、分散型オペレーティング・システムにおける様々なサーバについて述べた。また、オブジェクトの堆積に適した遠隔手続き呼出しのインタフェース記述言語や、堆積可能サーバの構造についても述べた。こうして、オブジェクトの堆積の特徴と有効性を明らかにした。

オブジェクトの堆積は、選択可能化を直接的に実現している。応用プログラムは、堆積するオブジェクトを選択することで、必要な機能を選択することができる。また、外部サーバの存在は、オペレーティング・システムの応用プログラム化が行われていることを示している。

本研究では、新しいオペレーティング・システムの全体の機能の提案を行った。そして、それらの機能の部分的な実現を行った。具体的には、軽量プロセスを実現するライブラリ、仮想プロセッサを提供するカーネル、並列シェルのプロトタイプ、いく

つかの堆積可能サーバの実現を行った。

選択可能化、隔離、および、オペレーティング・システムの応用プログラム化は、マッピングとスケジューリングにおける2つの調和を実現する上で有効な手段であることが示された。しかしながら、本研究において、全ての問題が解決されたわけではない。残された問題としては、スケジューリングに関しては、3. 13節において述べた実時間性の保証があげられる。マッピングに関しては、5. 6節において述べた複数の分散型堆積可能オブジェクトのマッピング・コントローラの間での協調があげられる。これらの問題を解決するには、隔離の度合いを下げる必要があると思われる。

本研究では、主にプロセッサ資源に関して、競合を解消し調和を実現する方法を提案した。今後の研究の方向としては、プロセッサ以外の計算機資源に関して、競合を解消し調和を実現する方法を研究することがあげられる。たとえば、ディスク資源やメモリ資源があげられる。これらの資源については、並列／分散応用プログラムが現れる以前から、データベース応用プログラムの効率的な処理を実現する局面において個別に対応がなされてきた。今後は、複数の利用者や応用プログラムを扱う局面において、統合的な対応が可能な方法を開発することが望まれる。

ネットワーク資源に関する調和の研究も興味深い。応用固有の情報を利用することにより、効率的なネットワーク通信プリミティブを開発することが可能となる。一方、ネットワーク通信プリミティブは、カーネルにおいて実現する方が効率がよいが、信頼できない応用プログラムをカーネルに置くことは難しい。このような競合状態を解消し、調和を図る方法も、重要な研究課題となるであろう。

軽量プロセスに関しては、ネットワーク機能との統合に関する研究が考えられる。たとえば、マッピング・コントローラにおいて、軽量プロセスの移動の技術を用いて、軽量プロセスを単位として負荷均衡を実現する方法が考えられる。

オブジェクトの堆積では、クライアントとサーバが非対象になっていた。すなわち、クライアントは、オブジェクト識別子を保持し、サーバを知っているが、サーバは、クライアントの情報を一切保持していない。オブジェクトの堆積の変形として、クライアントとサーバの区別がなくなり、オブジェクトが対等に結合される方法が考えられる。これを、オブジェクトの連結 (object-cascading) と呼ぶことにする。オブジェクトの連結は、局面に応じてクライアントとサーバの立場が入れ替わるようなウィンドウ・システムで有効であると思われる。

今後、オペレーティング・システムは、並列／分散応用プログラムを始めとして、新たに登場する様々な種類の応用プログラムの要求に対応する必要がある。オペレーティング・システムを必要とする並列／分散ハードウェアの進歩も著しい。このような応用プログラムの要求とハードウェアの変化に対応するために、オペレーティング・システムの研究は、ますます重要になってくるとと思われる。

オペレーティング・システムにおいて、計算機資源の有効利用を図るための資源割当ての最適化は、旧来から最も重要な仕事の1つであった。しかしながら、並列／分



分散アプリケーションが登場した現在、オペレーティング・システムは、もはや資源割当ての最適化を行う唯一のプログラムではない。今後、資源割当てに関して、オペレーティング・システムとそれ以外の並列／分散アプリケーションの間で相互利用が可能な技術がますます増えてくることが予想される。さらに、アプリケーションとその実行環境であるオペレーティング・システムの動的な相互作用が重要になってくるであろう。本研究が、このような技術の相互利用と動的な相互作用の研究に貢献することを確信する。

## 謝辞

本論文は、筆者が筑波大学大学院博士課程工学研究科において行った研究をまとめたものである。筑波大学電子・情報工学系清木康助教授には、研究の開始から現在に至るまで御指導と御助言をして頂いた。同助教授が進めて来られたデータベースの並列処理システムは、単一プログラミング・システムとしての並列／分散応用プログラムの例として、本研究の契機となり、その開発を通じて、本研究への数多くの貴重な示唆を受けることができた。ここに、心から感謝の意を表す。東京大学理学部益田隆司教授には、様々な角度から御指導、御助言を頂いた。そして、オペレーティング・システムの魅力と大域的な思考方法を伝授していただいた。心から感謝の意を表す。筑波大学電子・情報工学系中田育男教授には、本研究に対する有益なご意見を頂いた。同学系佐々政孝教授（現在東京工業大学）には、研究発表に関する御指導をして頂き、また、同教授の研究成果であるコンパイラ生成系 R i e を使わせて頂いた。同学系井田哲雄教授には、形式化の重要性を教えて頂いた。同学系板野肯三助教授には、本研究に関して有益なコメントを頂き、また、研究発表に関する指導をして頂いた。以上の方々に心から感謝の意を表す。

本論文をまとめるに当たり、筑波大学構造工学系星野力教授、同電子・情報工学系亀田壽夫教授、並びに、電子技術総合研究所情報アーキテクチャ部弓場敏嗣部長には、有益なコメントを頂いた。心から感謝の意を表す。

日本 I B M の田胡和哉博士、並びに、東京大学理学部加藤和彦博士をはじめとする益田研究室の皆様には、様々な機会を通じて熱心に討論していただいた。心から感謝の意を表す。本研究は、筑波大学オペレーティング・システム研究室において、行われたものである。同研究室の卒業生である高野陽介博士、劉澎博士、福田宗弘氏、黒沢貴弘氏、関安宏氏、波内みさ氏をはじめとする多くの人々から協力を得た。また、筑波大学工学研究科佐藤聡氏、同情報学類苅部朋幸氏には、本研究の成果を利用して頂き、また、有益な議論を交わすことができた。以上の方々を始め、同研究室の学生の方々に感謝する。また、ソフトウェア研究室ドメインの計算機利用者の仲間であった中田研究室、井田研究室の学生の方々、並びに、筑波大学工学研究科西山博泰氏、かん暁微氏を始めとする板野研究室の学生の方々に感謝する。

本研究では、f j / J U N E T におけるネットワーク・ニュースと電子メールによる議論が不可欠であった。このような環境を整えて頂いた筑波大学電子・情報工学系大田友一教授を始めとするポストマスタ、ならびに、歴代のニュース・システム管理者の方々に深く感謝する。そして、箕原さん@慶応大学、久野さん@大塚、筑波大学、平野さん@電総研、河野さん@CSL, S o n y、桜井さん@富士通、住元さん@富士通研、藤田さん@オムロン、乾さん@オムロン、門林さん@大阪大学、山井さん@奈良高専をはじめとする f j / J U N E T の皆さんには、ネットワーク・ニュースや電子メールを通じて、熱心に討論に参加して頂いた。深く感謝する。

## 参考文献

- [ 1 ] V. Abrossimov, M. Rozier and M. Shapiro: "Generic Virtual Memory Management for Operating System Kernels", SOSPI2, ACM Operating System Review, Vol. 23, No. 5, pp. 123-136 (1989).
- [ 2 ] M. J. Accetta, R. V. Baron, W. Bolosky, D. B. Golub, R. F. Rashid, A. Tevanian, and M. V. Young: "Mach: A New Kernel Foundation for UNIX Development", Proceedings of USENIX 1986 Summer Conference, pp. 93-112 (1986).
- [ 3 ] T. Anderson, B. Bershad, E. Lazowska and H. Levy: "Scheduler Activations: Effective Kernel Support for the User-Level Management of Parallelism", SOSPI3, ACM Operating System Review, Vol. 25, No. 5, pp. 95-109 (1991).
- [ 4 ] M. J. Bach: "The Design of the UNIX Operating System", Prentice-Hall (1986).
- [ 5 ] H. E. Bal, J. G. Steiner, and A. S. Tanenbaum: "Programming Languages for Distributed Computing Systems", ACM Computing Surveys, Vol. 21, No. 3 (1989).
- [ 6 ] "Balance 8000 Parallel Programming", Sequent Computer Systems, Inc. (1985).
- [ 7 ] A. P. Black: "Supporting Distributed Applications: Experience with Eden", Proceedings of the Tenth Symposium on Operating Systems Principles, pp. 181-193 (1985).
- [ 8 ] D. Black: "Scheduling Support for Concurrency and Parallelism in the Mach Operating System", IEEE Computer, Vol. 23, No. 5, pp. 35-43 (1990).
- [ 9 ] L. Cardelli and P. Wegner: "On Understanding Types, Data Abstraction, and Polymorphism", ACM Computing Surveys, Vol. 17, No. 4, pp. 471-522 (1985).
- [ 1 0 ] N. Carriero and D. Gelernter: "How to Write Parallel Programs: A Guide to the Perplexed", ACM Computing Surveys, Vol. 21, No. 3, pp. 323-357 (1989).
- [ 1 1 ] J. S. Chase, F. G. Amador, E. D. Lazowska, H. M. Levy, and R. J. Littlefield: "The Amber System: Parallel Programming on a Network of Multiprocessors", SOSPI2, ACM Operating System Review, Vol. 23, No. 5, pp. 147-158 (1989).
- [ 1 2 ] D. R. Cheriton: "The V Distributed System", Communications of the ACM, Vol. 31, No. 3, pp. 314-333 (1988).
- [ 1 3 ] R. Chin and S. Chanson: "Distributed Object-Based Programming Systems", ACM Computing Surveys, Vol. 23, No. 1, pp. 91-124 (1991).
- [ 1 4 ] S. Danforth and C. Tomlinson: "Type Theories and Object-Oriented Programming", ACM Computing Surveys, Vol. 20, No. 1, pp. 29-72 (1988).
- [ 1 5 ] P. Dasgupta, R. J. LeBlanc, Jr. and W. F. Appelbe: "The Clouds Distributed Operating System: Functional Description, Implementation Details and Related Work", Proc. IEEE 8th Intl. Conf. on Distributed Computing Systems, pp. 2-9 (1991).

- [ 1 6 ] J. D. Day: "The OSI Reference Model", Proc. of the IEEE, Vol. 71, No. 7, pp. 1334-1340 (1983).
- [ 1 7 ] R. Draves, B. Bershad, R. Rashid, and R. Dean: "Using Continuations to Implement Thread Management and Communication in Operating Systems", SOSPI3, ACM Operating System Review, Vol. 25, No. 5, pp. 122-137 (1991).
- [ 1 8 ] M. A. Ellis and B. Stroustrup: "The Annotated C++ Reference Manual", Addison-Wesley (1990).
- [ 1 9 ] C. J. Fleckenstein and D. Hemmedinger: "Using a Global Name Space for Parallel Execution of UNIX Tools", Communications of the ACM, Vol. 32, No. 9, pp. 1085-1090 (1989).
- [ 2 0 ] "Forth-83 Standard", Forth Standards Team (1983).
- [ 2 1 ] A. Fukuda, K. Murakami and S. Tomita: "Toward Advanced Parallel Processing: Exploiting Parallelisms at Task and Instructuon Levels", IEEE Micro, Vol. 11, No. 8, pp. 16-19, 50-61 (1991).
- [ 2 2 ] A. Goscinski: "Distributed Operating Systems -- The logical Design", Addison-Wesley, (1991).
- [ 2 3 ] R. Govindan and D. Anderson: "Scheduling and IPC Mechanisms for Continuous Media", SOSPI3, ACM Operating System Review, Vol. 25, No. 5, pp. 68-80 (1991).
- [ 2 4 ] R. G. Guy, J. S. Heidemann, Wai Mak, T. W. Page, Jr., G. J. Popek and D. Rothmeier: "Implementation of the Ficus replicated file system", USENIX 1990 Summer Conf., pp. 63-71 (1990).
- [ 2 5 ] S. Habert, L. Mosseri, and V. Abrossimov: "COOL: Kernel Support for Object-Oriented Environments", Proc. OOPSLA/ECOOP'90 Conf., ACM SIGPLAN Notes, Vol. 25, No. 10 (1990).
- [ 2 6 ] 萩野, 山岸: "T o M マイクロカーネル", 電子情報通信学会論文誌 D - 1, Vol. J75-D-1, No. 8, pp. 555-562 (1992).
- [ 2 7 ] 箱根, 横山, 谷口: "軽量プロセスと非完了システムコール機能の比較: 並列処理の性能評価", 情報処理学会第 4 2 回全国大会講演論文集 ( 4 ), 4K-11, pp. 15-16 (1991).
- [ 2 8 ] B. Hansen: "The Architecture of Concurrent Programs", Prentice-Hall (1997).
- [ 2 9 ] C. A. R. Hoare: "Communicating Sequential Processes", Prentice-Hall (1985).
- [ 3 0 ] S. Inohara, K. Kato, A. Narita and T. Masuda: "A Thread Facility Based on User/Kernel Cooperation in the XERO Operating System", Proc. 15th Intl. Computer Software & Applications Conf., pp. 398-405 (1991).
- [ 3 1 ] 井上外志雄: "標準 F O R T H", 共立出版 (1985).
- [ 3 2 ] 乾: "分散 O S M a c h とそのインプリメンテーション", 情報処理学会研究会報告, 90-OS-49-1 (1990).

- [ 3 3 ] M. B. Jones: "Inheritance in Unlikely Places: Using Objects to Build Derived Implementations of Flat Interfaces", Proc. 2nd International Workshop on Object Orientation in Operating Systems, pp.341-345 (1992).
- [ 3 4 ] M. B. Jones : "A Toolkit for Interposing User Code at the System Interface", Workshop on Operating Systems and Object Orientation at OOPSLA/EC00P'90 (1990).
- [ 3 5 ] A. R. Karlin, Kai Li, M. S. Manasse and S. Owicki: "Empirical Studies of Competitive Spining for a Shared-Memory Multiprocessor", SOSPI3, ACM Operating System Review, Vol.25, No. 5, pp.41-55 (1991).
- [ 3 6 ] 片上, 新城, 清木: "並列処理環境に適したプロセス・サーバと並列シユルの実現", 情報処理学会第44回全国大会講演論文集, 2G-2, pp.17-18 (1992).
- [ 3 7 ] J. Kepecs and M. Solomon: "SODA: A Simplified Operaing System for Distributed Applications", ACM Operating System Review, Vol.19, No.4, pp.45-56 (1985); originally presented at the Third ACM SIGACT/SIGOPS Symp. on Principles of Distributed Computing (1984).
- [ 3 8 ] Y. Kiyoki, T. Kurosawa, K. Kato, and T. Masuda: "The Software Architecture of a Parallel Processing System for Advanced Database Applications", Proc. 7th IEEE Conf. on Data Engineering, pp.220-229 (1991).
- [ 3 9 ] Y. Kiyoki, K. Kato and T. Masuda: "A Relational Database Machine Based on Functional Programming Concepts", Proc. ACM-IEEE Computer Society Fall Joint Computer Conf., pp.969-978 (1986).
- [ 4 0 ] Y. Kiyoki, K. Kato, N. Yamaguchi and T. Masuda: "A Stream-Oriented Approach to Parallel Processing for Deductive Database, Proc. 5th Intl. Workshop on Database Machines, pp.102-115 (1987).
- [ 4 1 ] 清木, 加藤: "関数型計算モデルのデータベース処理への適用", 情報処理, Vol.29, No. 8, pp.897-907(1988).
- [ 4 2 ] 久野: "型階層をもつオブジェクト指向言語", 情報処理, Vol.29, No.4, pp.318-324 (1988).
- [ 4 3 ] 前川, 所, 清水(編): "分散型オペレーティングシステム", 共立出版 (1991).
- [ 4 4 ] B. Meyer: "Object-oriented Software Construction", Interactive Software Engineering (1988).
- [ 4 5 ] B. W. Lampson: "Designing a Global Name Service", Proc. Symp. on Principles of Distributed Computing, pp.1-10 (1986).
- [ 4 6 ] S. J. Leffle, M. K. McKusic, M. J. Karels and J. S. Quarterman: "The Design and Implementation of the 4.3BSD UNIX Operating System", Addison-Wesley (1989).
- [ 4 7 ] E. Levy and A. Silberschatz: "Distributed File Systems: Concepts and Examples", ACM Computing Surveys, Vol.22, No. 4, pp.321-374 (1990).
- [ 4 8 ] H. Lieberman: "Using Prototypical Objects to Implement Shared Behavior in Object Oritented Systems", OOPSLA'86, pp.214-223 (1986).

- [ 4 9 ] B.Liskov: "Distributed Programming in ARGUS", Communications of the ACM, Vol.31, No.3, pp.300-312 (1988).
- [ 5 0 ] P.Liu, Y.Kiyoki and T.Masuda: "Efficient Algorithms for Resource Allocation in Parallel and Distributed Query Processing Environments", Proc. IEEE Intl. Conf. on Distributed Computing Systems (1989).
- [ 5 1 ] J.Magee, J.Kramer and M.Sloman: "Constructing Distributed Systems in Conic", IEEE Transactions on Software Engineering, Vol.15, No.6, pp.663-675 (1989).
- [ 5 2 ] "MC88200 Cache/Memory Management Unit User's Manual", Motorola (1990).
- [ 5 3 ] B.Marsh and M.Scott: "First-Class User-Level Threads", SOSP13, ACM Operating System Review, Vol.25, No.5, pp.110-121 (1991).
- [ 5 4 ] S.Mullender, G.Rossum, A.Tanenbaum, R.Renesse and H.Staveren: "Amoeba: A Distributed Operating System for the 1990s", IEEE Computer, Vol.23, No.5, pp.44-53 (1990).
- [ 5 5 ] S.J.Mullender and A.S.Tanenbaum: "The Design of a Capability-Based Distributed Operating System", The Computer Journal, Vol.29, No.4, pp.289-299 (1986).
- [ 5 6 ] T.Nakajima and H.Tokuda: "Implementation of Scheduling Policies in Real-Time Mach", Proc. 2nd International Workshop on Object Orientation in Operating Systems, pp.165-169 (1992).
- [ 5 7 ] 波内, 清木, 劉: "演繹データベースの並列処理方式の実現と資源割当て方式", 情報処理学会データベース・システム研究会, 89-DBS-72, pp.105-112, July 20-21 (1989).
- [ 5 8 ] "The NeXT System Reference Manual", NeXT, Inc. (1988).
- [ 5 9 ] "Network Programming", Sun Microsystems, Inc. (1990).
- [ 6 0 ] J.K.Ousterhout, A.R.Chersonson, F.Douglis, M.N.Nelson and B.B.Welch: "The Sprite Network Operating System", IEEE Computer, Vol.21, No.2, pp.23-36 (1988).
- [ 6 1 ] T.W.Page Jr., G.J.Popek and R.G.Guy: "Stackable Layers: An Object-Oriented Approach to Distributed File System Architecture", Workshop on Operating Systems and Object Orientation at OOPSLA/ECOOP (1990).
- [ 6 2 ] D.A.Patterson, G.Gibson and R.H.Katz: "A Case for Redundant Arrays of Inexpensive Disks (RAID)", Proc. of ACM SIGMOD Conf., pp.109-116 (1988).
- [ 6 3 ] R.Pike, D.Presotto, K.Thompson and H.Trickey: "Plan 9 from Bell Labs", proceedings of UKUUG, (1990).
- [ 6 4 ] G.J.Popek, B.Walker, J.Chow, D.Edwards, C.Kline, G.Rudisin and G.Thiel: "LOCUS: A Network Transparent, High Reliability Distributed System", Proceedings of the 8th Symposium on Operating Systems Principles, pp.160-168 (1981).
- [ 6 5 ] G.J.Popek: "The LOCUS Distributed System Architecture", The MIT Press, Cambridge USA (1986).

- [ 6 6 ] R. Rashid, A. Tevanian, Jr., M. Young, D. Golub, R. Baron, D. Black, W. J. Bolosky, and J. Chew: "Machine-Independent Virtual Memory Management for Paged Uniprocessor and Multiprocessor Architectures", IEEE Transactions on Computers, Vol. 37, No. 8, pp. 896-908 (1988).
- [ 6 7 ] R. Rashid: "Threads of a New System", UNIX REVIEW, No. 8, pp. 37-49 (1986).
- [ 6 8 ] "Reference Manual for the Ada Programming Language", United States Department of Defense (1983).
- [ 6 9 ] "Remote Procedure Call Programming Guide", Sun Microsystems, Inc. (1988).
- [ 7 0 ] D. S. H. Rosenthal: "Evolving the Vnode Interface", USENIX 1990 Summer Conf. (1990).
- [ 7 1 ] 酒井: "S P A R C プロセッサにおける高速コンテキスト切替え方式", 情報処理学会第4回コンピュータシステムシンポジウム論文集, 情報処理学会シンポジウム論文集, Vol. 92, No. 7, pp. 85-92 (1992).
- [ 7 2 ] 佐藤, 清木: "関数型計算に基づく並列型データベースシステムの並列処理実行系の実現方式", 情報処理学会データベース・システム研究会, 92-DBS-89-19, July 24 (1992).
- [ 7 3 ] 関, 関島, 清木: "並列処理システム S M A S H における機械系 C A D データベース実現のための一考察", 情報処理学会データベース・システム研究会, 89-DBS-72, pp. 177-184, July 20-21 (1989).
- [ 7 4 ] L. Sha, R. Rajkumar and J. Lehoczky: "Priority Inheritance Protocols: An Approach to Real-Time Synchronization", IEEE Transactions on Computers, Vol. 39, No. 9, pp. 1175-1185 (1990).
- [ 7 5 ] K. Shimizu, M. Maekawa and J. Hamano: "Hierarchical Object Groups in Distributed Operating Systems", Proc. IEEE 8th Intl. Conf. on Distributed Computing Systems, pp. 18-24 (1991).
- [ 7 6 ] Y. Shinjo and Y. Kiyoki: "Harmonization methods in the ReSC Distributed Operating System", Journal of Concurrency: Practice and Experience, special issue on HPDC-1 (投稿中).
- [ 7 7 ] 新城, 清木: "仮想プロセッサを提供するオペレーティング・システム・カーネルの構成法", 情報処理学会論文誌, Vol. 34, No. 3 (1993).
- [ 7 8 ] Y. Shinjo and Y. Kiyoki: "The Object-Stacking Model for Structuring Object-Based Systems", Proc. 2nd International Workshop on Object Orientation in Operating Systems (I-WOOS' 92), pp. 328-340 (1992).
- [ 7 9 ] Y. Shinjo and Y. Kiyoki: "Harmonizing a Distributed Operating System with Parallel and Distributed Applications", Proc. 1st International Symposium on High Performance Distributed Computing (HPDC-1), pp. 114-123 (1992).
- [ 8 0 ] 新城, 清木: "オブジェクトの堆積に適した R P C スタブ生成器", 情報処理学会第44回全国大会講演論文集(4), 2G-1, pp. 15-16 (1992年)
- [ 8 1 ] 新城, 清木: "並列プログラムを対象とした軽量プロセス実現方式", 情報処理学会論文誌, Vol. 33, No. 1, pp. 64-73 (1992).

- [ 8 2 ] Y. Shinjo and Y. Kiyoki: "ReSC: A Distributed Operating System for Parallel and Distributed Applications", Proc. First International Conference on Parallel and Distributed Information Systems (PDIS), p. 171 (1991).
- [ 8 3 ] 新城, 清木: "分散型オペレーティング・システムにおけるオブジェクトの堆積", 情報処理学会第42回全国大会講演論文集(4), 3K-8, pp.15-16 (1991).
- [ 8 4 ] 新城: "分散型OS ReSCにおけるデータベース処理の支援", 電子情報通信学会データベースワークショップ論文集, pp.118-123 (1991)
- [ 8 5 ] Y. Shinjo and Y. Kiyoki: "Multiple Views of ReSC Distributed Operating System", Workshop on Operating Systems and Object Orientation at OOPSLA/ECOOP'90 (1990).
- [ 8 6 ] 新城, 清木, 益田: "並列・分散処理環境を対象としたOS ReSCにおける分散最適化", 情報処理学会第40回全国大会講演論文集, 5G-3 (1990).
- [ 8 7 ] 新城, 清木: "ReSC: 並列アプリケーションのためのオペレーティング・システム", 情報処理学会第39回全国大会講演論文集, 6P-2, pp.1247-1248 (1989).
- [ 8 8 ] 新城, 清木: "データベースの並列処理を支援するオペレーティング・システムの基本機能", 情報処理学会研究会報告, 89-OS-44, 89-DBS-73 (1989).
- [ 8 9 ] 新城, 清木, 劉, 益田: "データベースおよび知識ベースを対象としたストリーム指向型並列処理系の共有メモリ・マシン上への実現", 情報処理学会アドバンストデータベースシステム・シンポジウム論文集, Vol.88, No.9, pp.117-126 (1988).
- [ 9 0 ] 新城: "密結合並列処理システム上でストリーム指向型関係演算処理をする際の資源割当て", 第21回情報処理学会若手の会シンポジウム, 湯河原 (1988).
- [ 9 1 ] 新城, 田胡, 高野, 福田, 関: "分散システムAgora-Iのユーザ環境", 情報処理学会第36回全国大会講演論文集, 4D-9, pp.325-326 (1988).
- [ 9 2 ] M. Sloman and J. Kramer: "Distributed Systems and Computer Networks", Prentice-Hall International(UK) Ltd. (1987).
- [ 9 3 ] "SunOS Reference Manual", Sun Microsystems, Inc. (1988).
- [ 9 4 ] V. Srinivasan and J.C. Mogul: "Spritely NFS: Experiments with Cache-Consistency Protocols", SOSPI2, ACM Operating System Review, Vol.23, No.5, pp.45-57 (1989).
- [ 9 5 ] M. Stonebraker: "Operating System Support for Database Management", Communications of the ACM, Vol.24, No.7, (1981).
- [ 9 6 ] 田胡, 益田: "分散型オペレーティング・システム", 情報処理, Vol.28, No.4, pp.437-445 (1987).
- [ 9 7 ] 田胡, 益田: "オペレーティング・システムの構造記述に関する一試み", 情報処理学会論文誌, Vol.25, No.4, pp.524-534 (1984).



- [ 9 8 ] A. S. Tanenbaum and R. van Renesse: "Distributed Operating Systems",  
ACM Computing Surveys, Vol.17, No.4, pp.419-470 (1985).
- [ 9 9 ] A. S. Tanenbaum: "Modern Operating Systems", Prentice-Hall (1992).
- [ 1 0 0 ] C. P. Thacker, L. C. Stewar, E. H. Satterhwaite, Jr.: "Firefly: A  
Multiprocessor Workstation", IEEE Transaction on Computers, Vol.37,  
No.8, pp.909-920(1988).
- [ 1 0 1 ] "Transmission Control Protocol", DARPA Internet Program Protocol  
Specification, RFC 793 (1981).
- [ 1 0 2 ] H. Tokuda and T. Nakajima: "Evaluation of Real-Time Synchronization  
in Real-Time Mach", USENIX 2nd Mach Symposium (1991).
- [ 1 0 3 ] A. Tucker and A. Gupta: "Process Control and Scheduling Issues for  
Multiprogrammed Shared-Memory Multiprocessors", SOSPI2, ACM  
Operating System Review, Vol.23, No.5, pp.159-166 (1989).
- [ 1 0 4 ] 梅村, 大里 : "L i s p上のオブジェクト指向プログラミング", 情報処  
理, Vol.29, No.4, pp.303-309 (1988).
- [ 1 0 5 ] "UniOS-Mach Reference Manual", オムロン (1990).
- [ 1 0 6 ] M. Weiser, A. Demers and C. Hauser: "The Portable Common Runtime  
Approach to Interoperability", SOSPI2, ACM Operating System Review,  
Vol.23, No.5, pp.114-122 (1989).
- [ 1 0 7 ] "Workshop on Object-Orientation in Operating Systems",  
OOPSLA/ECOOP'90 Addendum to the Proceedings, edited by V. Russo and  
M. Shapiro, pp.81-91 (1990).
- [ 1 0 8 ] W. Wulf, E. Cohen, W. Corwin, A. Jones, R. Levin, C. Pierson, and  
F. Pollack: "HYDRA: The Kernel of a Multiprocessor Operating System",  
Communications of the ACM, Vol.17, No.6, pp.337-345 (1974).
- [ 1 0 9 ] Y. Yokote, F. Teraoka, A. Mitsuzawa, N. Fujinami and M. Tokoro:  
"The Muse Object Architecture: A New Operating System Structuring  
Concept", ACM Operating System Review, Vol.25, No.2, pp.22-46  
(1991).
- [ 1 1 0 ] C. T. Yu and C. C. Chang: "Distributed Query Processing", ACM  
Computing Surveys, Vol.16, No.4, pp.399-433, (1984).

# 付録の目次

付録A	遠隔手続き呼出し	1
A-1	オブジェクト識別子	1
A-2	マスク	1
A-3	利用者識別子	1
A-4	アクセス権のチェックの例	2
A-5	カーネル・コール	3
A-6	サーバ・プログラムの例	4
A-6.1	手続きの表	4
A-6.2	セッション構造体	5
A-6.3	遠隔手続き呼出し受け付け軽量プロセス	5
A-6.4	遠隔手続き呼出し入口手続き	6
A-6.5	メモリ中のオブジェクトに対する手続き	6
A-7	クライアント・プログラムの例	7
付録B	ReSCシステムの標準的なオブジェクトとサーバ	9
B-1	インタフェース記述言語proc	9
B-1.1	構文	9
B-1.2	基本データ型	10
B-1.3	一般的な型構成子	10
B-1.4	拡張した型構成子	11
B-2	オブジェクト手続きやサーバ手続きで利用されるデータ型	12
B-3	一般のサーバ	13
B-4	一般のオブジェクト	14
B-5	ファイル・サーバとファイル・オブジェクト	15
B-6	ディレクトリ・サーバとディレクトリ・オブジェクト	16
B-7	メモリ・サーバとメモリ・オブジェクト	17
B-8	プロセス・サーバとプロセス・オブジェクト	18
B-9	カーネル・サーバ群	19

付録C	マイクロプロセス・ライブラリの機能	21
C-1	全体の構成	21
C-2	データ型	21
C-3	第0層 カーネル・コール	22
C-4	第1層 スピンロックとコンテキスト切替え	23
C-4.1	スピンロック	23
C-4.2	コンテキストの生成と切替え	23
C-5	第2層 マイクロプロセスの生成と基本的な同期	23
C-5.1	現在実行中のマイクロプロセス	24
C-5.2	マイクロプロセスの生成・消滅	24
C-5.3	マイクロプロセス間の同期	25
C-6	第3層 高度な同期通信プリミティブ	25
C-6.1	セマフォ	25
C-6.2	モニタ	26
C-6.3	ランデブ	26
C-6.4	ストリーム	27
C-7	マイクロプロセス・スケジューラ	28
C-8	レディ・キュー・モジュール	28
付録D	プロセス操作の性能を計測するためのベンチマーク	31
D-1	プロセスの生成・消滅	31
D-2	コンテキスト切替えを伴わないプロセス間の同期	32
D-3	コンテキスト切替えを伴うプロセス間の同期	32
D-3.1	セマフォを用いたプログラム	32
D-3.2	モニタと条件変数を用いたプログラム	33
D-3.3	コンテキスト切替えを起こすプリミティブを用いたプログラム	34

## 付録 A 遠隔手続き呼出し

ここでは、遠隔手続き呼出し関連のカーネル・コールと、それを利用する典型的なサーバとクライアントの構造を、C言語を用いて示す。

### A-1 オブジェクト識別子

システム・レベルのオブジェクトは、全て一様なオブジェクト識別子により識別される。本システムのオブジェクト識別子は、以下の `oid_t` で示された構造を持つ。

```
typedef int    site_t ;        /* サイト識別子 */

typedef int    server_t ;     /* サーバ識別子 */

typedef struct
{
    char[32];                /* 32 == NFS_FHSIZE */
} inum_t ;                  /* オブジェクト番号 */

typedef struct
{
    site_t      site ;
    server_t    server ;
    inum_t      inum ;
} oid_t ;                  /* オブジェクト識別子 */
```

### A-2 マスク

本システムでは、様々な場所でマスク (mask) を用いる。マスクとは、対応する整数の有効範囲を現すものである。マスクのあるビットが1であることは、そのビットが有効であることを意味する。

```
#define All1    (~0)    /* '~' は、ビット反転を行う演算子 */
```

全ビットが1のマスクは、対応する整数の全てのビットが有効であることを意味する。

### A-3 利用者識別子

利用者識別子は、オブジェクトへのアクセス権の確認、セキュリティのチェックにおいて、利用者を識別するものである。利用者識別子は、以下の `uid_t` に示された構造を持つ。

```
typedef struct
{
    int    u_no ;            /* 利用者番号 */
    int    u_mask ;        /* マスク */
} uid_t ;
```

利用者番号 `u_no` は、各利用者に一意になるように与えられた整数である。マスク `u_mask` は、利用者番号の有効範囲を示す。

本システムのプロセスは、属性として利用者識別子を持っている。利用者識別子の利用者番号の、マスクにより残された部分がアクセス権のチェックに用いられる。たとえば、システム管理者の場合、マスクとして0が用いられる。これにより、全てのアクセス権のチェックが無効になり、任意のオブジェクトにアクセスすることが可能になる。一般の利用者のマスクは、ALL1 である。

#### A-4 アクセス権のチェックの例

この節では、A-3 節で述べた利用者識別子を用いたアクセス権のチェックの例を示す。たとえば、次の様なオブジェクトを考える。

```
struct an_object
{
    int    owner ;
    int    rmask ;
    int    wmask ;
    int    xmask ;
    . . .
};
```

属性 rmask, wmask, xmask が、それぞれ読み込み、書き込み、実行のアクセス可能性を決定する変数である。ある利用者識別子 uid を持つクライアント・プロセスのアクセス権を考える。

```
struct process
{
    uid_t uid ;
    . . .
};
```

たとえば、読み込みの場合、次の式が真ならば、その利用者はオブジェクトへの読み込みのアクセス権を持っていると判定する。

```
(uid.u_no & uid.u_mask & rmask) ==
(owner    & uid.u_mask & rmask)
```

このように、クライアント・プロセスの利用者識別子のマスクとオブジェクトのマスクのビットが1の部分について、クライアント・プロセスの利用者番号 u\_no とオブジェクトの所有者の利用者番号 owner を比較して、アクセス権の有無を判定する。したがって、オブジェクトのマスクとして0を指定することにより、全ての利用者からのアクセスを許すことを意味する。ALL1 を指定することにより、オブジェクトの所有者にみのアクセスを許すことを意味する。0 と ALL1 の中間を指定することにより、あるグループによるアクセスを許すことを意味する。

例として、次の様なアクセスが可能なオブジェクトを生成ことを考える。

- 全ての利用者からの読み込みを許す。
- そのオブジェクトの所有者だけに書き込みを許す。
- 同じグループに属する利用者の実行を許す。

これは、UNIXにおいて、ファイルのモードが 754 ( rwxr-xr-- ) として実現される。ここでは、利用者番号とそのマスクを32ビットの整数、その中で、グループの管理のために、上位16ビットを用いることにする。上記のようなアクセスが可能

なオブジェクトを実現するためには、次のように、マスクを設定する。

```
rmask = 0 ;
wmask = ALL1 ;
xmask = ALL1 << 16 ;
```

#### A-5 カーネル・コール

R e S C システムの遠隔手続き呼出し機能は、次の5つのカーネル・コール `rpc_send()`, `rpc_receive()`, `rpc_reply()`, `rpc_forward()`, `rpc_send_uid()` により提供される。最初のものが、クライアントにより、残りの4つがサーバにより利用される。特に、最後の2つは、堆積可能なサーバにおいて利用される。

```
int rpc_send( site, server, procno, arg, res, timeout )
site_t site ;
server_t server ;
int procno ;
rpcbuff_t *arg ;
rpcbuff_t *res ;
time_t *timeout ;
```

`rpc_send()` は、クライアントが遠隔手続き呼出しの呼出し（要求の送信と結果の受信）を行うカーネル・コールである。`site` は、サーバが存在するサイト、`server` は、サーバ識別子（遠隔手続き呼出しのポート番号）である。オブジェクト手続きの場合、`site` と `server` は、オブジェクト識別子から抜き出されたものを用いる。`procno` は、遠隔手続き呼出しの手続き番号である。`arg` は、クライアントからサーバに渡される引数を格納するバッファを、次のような構造体により指定する。

```
typedef struct
{
    char      *m_addr ;
    int       m_len ;
} rpcbuff_t ;
```

`res` は、サーバからクライアントへ返される結果を格納するバッファを指定する。`timeout` は、遠隔手続き呼出しの完了の最大待ち時間を指定する。

```
typedef int rpcid_t ;
rpc_id rpc_receive( server, procno, arg, res, timeout, uid )
server_t server ;
int *procno ;
rpcbuff_t *arg ;
rpcbuff_t *res ;
time_t *timeout ;
uid_t *uid ;
```

`rpc_receive()` は、サーバ側において遠隔手続き呼出しの受け付けを行うカーネル・コールである。`server`, `arg`, `res`, `timeout` の意味は、`rpc_send()` と同じである。`uid` は、クライアント・プロセスの利用者識別子である。これはカーネルによって自動的に付加されるものである。サーバは、これを利用してアクセス権のチェックを行うことができる。`rpc_receive()` の結果は、遠隔手続き呼出しの識別子である。これは、次の `rpc_reply()`, `rpc_forward()` において使用される。

```
int rpc_reply( rpcid, errno )
rpcid_t rpcid ;
```

```
int errno ;
```

rpc\_reply() は、サーバ側において遠隔手続き呼出しの応答を行うカーネル・コールである。rpcid は、rpc\_receive() により返された識別子を指定する。errno は、呼出し側の rpc\_send() カーネル・コール自身の結果として、クライアントに渡される。errno が 0 であることは、エラーが発生しなかったことを示す。遠隔手続き呼出しの結果は、rpc\_reply() を行う前に、rpc\_receive() により受け付け時に res で指定された領域に格納されていなければならない。

```
int rpc_forward( rpcid, site, server, procno, timeout, uid )
```

```
rpcid_t rpcid ;  
site_t site ;  
server_t server ;  
int procno ;  
time_t *timeout ;  
uid_t *uid ;
```

rpc\_forward() は、遠隔手続き呼出しの転送を行うカーネル・コールである。rpcid に指定された遠隔手続き呼出しのセッションを、site, server をで指定されたサーバの手続き番号 procno の手続きに転送する。ここで、uid には、転送するとき用いるクライアントの利用者識別子である。これを書き換えることで、受け付けた時とは異なる権限で遠隔手続き呼出しの呼出しを行うことが可能となる。ただし、書き換えが可能な部分は、そのプロセスの利用者識別子のマスクが 0 の部分だけである。

```
int rpc_send_uid( site, server, procno, arg, res, timeout, uid )
```

```
site_t site ;  
server_t server ;  
int procno ;  
rpcbuff_t *arg ;  
rpcbuff_t *res ;  
time_t *timeout ;  
uid_t *uid ;
```

rpc\_send\_uid() は、rpc\_send() と同じく、遠隔手続き呼出しの呼出しを行うカーネル・コールである。rpc\_send() との違いは、uid を指定することができる点にある。uid は、rpc\_forward() と同じく、遠隔手続き呼出しの呼出しを行う時の利用者識別子である。

## A-6 サーバ・プログラムの例

以下にサーバ・プログラムの例を示す。重要なデータ構造としては、手続きの表、セッション構造体がある。重要な手続きとしては、遠隔手続き呼出し受け軽量プロセス、遠隔手続き呼出し入口手続き、メモリ中のオブジェクトに対する手続きがある。

### A-6.1 手続きの表

```
struct rpc_proc  
{  
    int    procno ;  
    int    (*xdr_arg)();  
    int    arg_size ;  
    int    (*xdr_res)();  
    int    res_size ;  
    char  *(*local)();  
};
```

```

typedef struct rpc_proc *rpc_proc_t ;

struct rpc_proc rpc_proctab[] =
{
    {SERVER_NULL, xdr_void, 0,
     xdr_void, 0, server_null },
    {SERVER_CREATE_WITH_LOWER, xdr_oid_array_t, sizeof(oid_array_t),
     xdr_oid_res_t, sizeof(oid_res_t), server_create_with_lower},
    . . .
};
int rpc_proctab_nentries = sizeof(rpc_proctab)/sizeof(rpc_proc) ;

extern rpc_proc_t      rpc_proctab_find( /*tab,nent,procno*/ );

```

rpc\_proctab は、そのサーバが受け付け可能な遠隔手続き呼出しの手続きの表である。procnoは、手続き番号、xdr\_arg は、引数を解く (unpack, unmarshall) 手続き、arg\_size は、引数の大きさ、xdr\_res は、結果をまとめる (pack, marshal) 手続き、res\_size は、結果の大きさ、local は、メモリ中のオブジェクトへの手続きのポインタである。rpc\_proctab\_find() は、表を検索するライブラリ関数である。

#### A-6.2 セッション構造体

```

struct rpc_session_t
{
    rpcid_t      rpcid ;
    rpcbuff_t    arg ;
    rpcbuff_t    res ;
    int          procno ;
    uid_t        uid ;
};

```

これは、遠隔手続き呼出しのセッションの情報を保持する構造体である。この構造体は、A-5 で述べた rpc\_receive() の引数と結果に対応している。

#### A-6.3 遠隔手続き呼出し受け付け軽量プロセス

```

rpc_acceptance_lwp()
{
    rpc_proc_t    p ;
    rpc_session_t *s ;

    while( 1 )
    {
        if( s == 0 )
            s = malloc( sizeof(rpc_session_t) );
        s->rpcid =
            rpc_receive( SERVER1, &s->procno, &s->arg, &s->res,
                        &TIMEOUT, &s->uid );
        p = rpc_proctab_find( rpc_proctab, rpc_proctab_nentries,
                             s->procno );
        if( p )
        {
            mp_create( rpc_entry_procedure, STACK_SIZE, M_READY, s, p );
            s = 0 ;
        }
    }
}

```



```

    }
    else
    {
        rpc_reply( s->rpcid, ERROR_NOPROC );
    }
}
}

```

これは、遠隔手続き呼出し受け付け軽量プロセスの骨格を示している。この軽量プロセスは、このように無限ループを実行しており、通常、`rpc_receive()` カーネル・コールを発行し、遠隔手続き呼出しの受け付けを行っている。クライアントから遠隔手続き呼出しを受け付けると、`rpc_proctab_find()` により、手続きの表を検索する。対応する手続きが見つかった場合、遠隔手続き呼出しエントリ手続き `rpc_entry_procedure()` をトップとする軽量プロセスを生成する。見つからなかった場合、エラーを返す。

#### A-6.4 遠隔手続き呼出し入口手続き

```

void rpc_entry_procedure( s, p )
    rpc_session_t *s ;
    rpc_proc_t    p ;
{
    char          *arg ;
    char          *res ;

    arg = arg_decode( p->xdr_arg, p->arg_size, s->arg );
    res = (*p->local)( arg, s, p );
    res_encode( p->xdr_res, res, s->res );
    rpc_reply( s->rpcid, 0 );
    arg_free( p->xdr_arg, arg );
    res_free( p->xdr_res, res );
    free( s );
    return; /* mp_exit() */
}

```

これは、遠隔手続き呼出し入口手続きである。この手続きの役割は、クライアントから送られてきた引数を解き、メモリ中のオブジェクトに対する手続きを呼び出し、結果をクライアントへ返すことである。まず、ライブラリ関数 `arg_decode()` により、遠隔手続き呼出しの引数を解く。続いて、対応する局所オブジェクトに対する手続き `p->local` を呼び出す。そして、ライブラリ関数 `res_encode()` により、遠隔手続き呼出しの結果をまとめ、`rpc_reply()` カーネル・コールにより、クライアントに結果を返す。その後、不用になった引数、結果、セッション構造体のメモリを開放する。最後に、この関数から復帰することで、現在の軽量プロセスを終了させている。

#### A-6.5 メモリ中のオブジェクトに対する手続き

```

object_forward( arg, s )
    some_arg_t    *arg ;
    rpc_session_t *s ;
{
    object_t      *obj ;

    obj = obj_fetch( arg );
    if( obj == 0 )
    {

```

```

        rpc_reply( s, ERROR_NOOBJECT );
        free( s );
        arg_free( arg );
        mp_exit();
    }

    obj_lock( obj );
    obj_access( obj, s->uid );
    arg->oid = obj->lower ;
    arg_encode( p->xdr_arg, arg, s->arg );
    rpc_forward( s->rpcid, obj->lower.site, obj->lower.server,
                s->procno, &TIMEOUT, &s->uid );
    obj_unlock( obj );
    obj_free( obj );
    arg_free( p->xdr_arg, arg );
    free( s );
    mp_exit();
}

```

これは、典型的な転送 (forward) を行うようなメモリ中のオブジェクトへの手続きである。まず、関数 `obj_fetch()` を呼び出し、メモリ中のオブジェクトを得る。失敗すると、エラーを返す。次に、オブジェクトにロックを掛け、アクセス権のチェックを行う。そして、引数の中のオブジェクト識別子の部分 (`arg->oid`) を、下位層のオブジェクト (`obj->lower`) に書き換え、`s->rpcid` で指定された遠隔手続き呼出しのセッションを、`obj->lower.site`, `obj->lower.server` で指定された、下位層のオブジェクトのサーバに転送 (`rpc_forward()`) している。その後、ロックの開放、オブジェクトの開放、セッション構造体の開放を行い、軽量プロセスを終了している。

#### A-7 クライアント・プログラムの例

以下に、クライアント・プログラムの例を示す。

```

getattr( oid, attr )
    oid_t *oid ;
    attr_t *attr ;
{
    memmaddr_t arg, res ;
    char arg_buff[MAX_ARG_BUFF_SIZE] ;
    char res_buff[MAX_RES_BUFF_SIZE] ;

    arg.m_addr = arg_buff ;
    arg.m_len = arg_encode( xdr_oid_t, oid, arg_buff );
    res.m_addr = res_buff ;
    res.m_len = MAX_RES_BUFF_SIZE ;

    errno = rpc_send( oid->site, oid->server, OBJECT_GETATTR,
                    &arg, &res, TIMEOUT );
    if( errno == 0 )
    {
        res_decode( xdr_attr_t, attr, arg_buff );
    }
    return( errno );
}

```

これは、オブジェクトの属性 (attr\_t) を得る手続きである。まず、引数 (oid) を、ライブラリ関数 arg\_encode() によりまとめ、arg\_buff に格納している。次に、結果を受け取るためのバッファ res\_buff を、res に設定している。そして、カーネル・コール rpc\_send() により、遠隔手続き呼出しを行っている。ここで、遠隔手続き呼出しにより呼び出すサーバを、oid->site, oid->server で、オブジェクト識別子の中から抜出している。OBJECT\_GETATTR は、手続き番号である。そして、遠隔手続き呼出しが成功すると (errno==0)、結果を解いている。

## 付録 B ReSCシステムの標準的なオブジェクトとサーバ

ここでは、5.8節で述べたインタフェース記述言語を用いて、ReSCシステムで標準的に利用されるオブジェクトとサーバのインタフェースについて述べる。ここでは、簡単のため、結果として返されるエラーを省略する。

### B-1 インタフェース記述言語 `proc`

5.8節で述べたインタフェース記述言語を `proc` と名付ける。これは、プログラミング言語 `fun` において関数を意味する `fun` と記述するところを、本インタフェース記述言語では、手続きを意味する `proc` と記述することに由来する。

#### B-1.1 構文

本言語におけるインタフェースは、次のような文の並びにより記述される。

<名前> : <型> = <値> ;

この構文は、`fun` に準じている。`fun` は、1階型付きλ計算に基づく、多相 (polymorphism) とオブジェクト指向言語をモデル化するために設計された言語である。ただし、`fun` において関数を意味する `fun` と記述するところを、本インタフェース記述言語では、手続きを意味する `proc`、あるいは、`sproc` と記述する。また、`fun` では、型と値が明確に区別されているのに対して、本インタフェース記述言語では、型を型型 (型を格納するための変数のための型) の値として扱う。たとえば、次の文：

```
oid_t:type =
struct
{
    site_id      : site_t;
    server_id    : server_t ;
    object_number : object_number_t ;
};
```

では、`oid_t` という変数が、型を格納するための変数 (`:type`) であり、`oid_t` の値として、`struct` 以下に記述されている構造体が束縛される。なお、この定義は、C言語における定義：

```
typedef struct {
    site_t      site_id;
    server_t    server_id;
    object_number_t object_number;
} oid_t ;
```

と同じ意味を持つ。

## B-1.2 基本データ型

本インタフェース記述言語において扱えるデータ型は、C言語、および、SunRPCのインタフェース記述言語（rpcgenコマンドの入力）を拡張したものになっている。基本データ型としては、次のようなものがある。

void型	void
整数型	int, unsigned, long
列挙型	enum { <名前>=値, ... }
実数型	float, double
文字型	char

## B-1.3 一般的な型構成子

C言語、および、SunRPCのインタフェース記述言語と同等の型構成子を以下に示す。

### ■配列

配列には、次のような種類がある。

array[n] of <型>	<型>型のn個の要素からなる固定長の配列
array<n> of <型>	<型>型の最大n個の要素からなる可変長の配列
array<> of <型>	<型>型の要素からなる可変長の配列 (最大要素数は不明)
array of <型>	array<> of <型> と同じ

### ■文字列型

文字列型は、文字列（0 終端された文字の配列）を表す型である。次のような種類がある。

string<n>	最大n文字の可変長の文字列
string<>	最大長が不明の可変長の文字列
string	string<> と同じ

### ■不定型

不定型（opaque）とは、内容を一切解釈しないバイト列である。次のような種類がある。

opaque[n]	nバイトの固定長のバイト列
opaque<n>	最大nバイトの可変長のバイト列
opaque<>	最大長が不明のバイト列
opaque	opaque<> と同じ

### ■構造体

これは、C言語の構造体に相当する。次のような構文により、構造体が定義される。

```
struct { <フィールド並び> }
```

ここで、<フィールド並び>とは、

```
<名前の並び> : <型> ;
```

の並びである。

## ■選択子付き共用体 (discriminated union)

これは、S u n R P Cのインタフェース記述言語における選択子付き共用体 (discriminated union) に対応するものである。これは、P a s c a l言語における可変長レコードに相当する。次のような構文により、選択子付き共用体が定義される。

```
union switch(<選択子>)
{
case <選択子1>:
    <フィールド並び1>;
case <選択子2>:
    <フィールド並び2>;
...
default:
    <デフォルトのフィールド並び>;
}
```

### B-1.4 拡張した型構成子

ここでは、C言語、および、S u n R P Cのインタフェース記述言語には存在しない型構成子 proc, sproc, set of, server について説明する。

#### ■オブジェクト手続き構成子 proc

オブジェクト手続き構成子 proc は、オブジェクトに対する遠隔手続き呼出しを定義する。次のような構文により、オブジェクト手続きの型が定義される。

```
proc( <引数のフィールド並び> )
    <結果のフィールド並び>
```

これは、以下で述べるサーバ手続きに展開される。

```
sproc( this:oid_t; <引数のフィールド並び> )
    <結果のフィールド並び>
```

このように、オブジェクト手続きは、第1引数として操作対象となるオブジェクトの識別子をとるようなサーバ手続きである。

#### ■サーバ手続き構成子 sproc

サーバ手続き構成子とは、サーバ手続き型の1つの型を作るものである。次のような構文により、サーバ手続きの型が定義される。

```
sproc( <引数のフィールド並び> ) <結果のフィールド並び>
```

<引数のフィールド並び>では、引数となる構造体を定義する。  
<結果のフィールド並び>では、結果となる構造体を定義する。この型の値は、遠隔手続き呼出しの手続き番号を表す整数である。たとえば、次の文：

```
copy: sproc( orig:oid_t ) new:oid_t = 2 ;
```

は、copy という名前の遠隔手続き呼出しの手続きを定義している。ここで、変数 copy の値として、2 が束縛されている。この2という値は、copy手続きの遠隔手続

き呼出しの手続き番号である。この手続きは、引数としてオブジェクト識別子を取り、結果としてオブジェクト識別子を返す。ここで、次の部分：

```
sproc( orig:oid_t ) new:oid_t
```

が、サーバ手続き型の値となる。したがって、上記の copy 手続きを次のように型変数を用いて宣言することも可能である。

```
proc_oid_oid_t: type = sproc( orig:oid_t ) new:oid_t ;  
copy: proc_oid_oid_t = 2 ;
```

なお、sproc により、遠隔手続き呼出しのクライアント側のスタブが生成される。

#### ■集合型 set of proc, set of sproc

集合型は、オブジェクト手続き、または、サーバ手続きの集合を作るための型である。この型の値は、以下で述べるサーバ型構成子の引数となる。次のような構文により、集合型が定義される。

```
set of {proc | sproc}
```

この型の値は、オブジェクト、または、サーバのインタフェースを表すために用いられる。たとえば、次の文：

```
any_server: set of sproc = { copy, create };
```

では、any\_server という名前の変数に、copy, create という要素からなる集合が束縛される。これは、任意のサーバ any\_server は、サーバのインタフェースとして、手続き copy と create を持つことを意味する。

#### ■サーバ型構成子 server

サーバ型構成子は、サーバのスタブを生成するための型構成子である。次のような構文により、サーバ型が定義される。

```
server( <サーバ手続きの集合> )
```

このように、<サーバ手続きの集合>を引数として取り、それらの手続きを受け付けるようなサーバのスタブを生成する。この型の変数の値は、サーバ識別子である。たとえば、次の文：

```
StdFS: server( file_server ) = 3 ;
```

では、StdFSは、サーバ型である。その値3は、サーバStdFSのサーバ識別子となる。そして、集合 file\_server に保持されているサーバ手続きを受け付けるサーバのスタブが生成される。

## B-2 オブジェクト手続きやサーバ手続きで利用されるデータ型

ここでは、付録B-1で述べたインタフェース記述言語procの基本データ型以外に、オブジェクト手続きやサーバ手続きで利用されるデータ型について述べる。

## ■オブジェクト識別子

```
site_t: type = int ;
server_t: type = int ;
object_number_t: type = opaque[32] ;

oid_t: type =
struct
{
    site_id      : site_t;
    server_id    : server_t ;
    object_number : object_number_t ;
};
```

これは、2.6節において導入したR e S Cのオブジェクト識別子を表している。site\_t は、サイト識別子を、server\_t は、サーバ識別子を、object\_number\_t は、オブジェクト番号を表す型である。そして、oid\_t は、オブジェクト識別子を意味する型である。現在、オブジェクト番号は、32バイトのバイト列である。そして、オブジェクト識別子全体の大きさは、40バイトである。

## ■オブジェクト識別子の配列

```
oid_array_t: type = array<> of oid_t ;
```

これは、オブジェクト識別子の可変長の配列である。

## ■ディレクトリ・エントリ

```
dir_entry_t: type = struct { oid:oid_t; name:string };
```

これは、ディレクトリの1つのエントリを示す型である。エントリとは、ここで示されているように、オブジェクト識別子と文字列の名前の組である。

## B-3 一般のサーバ

この節では、一般的なサーバが全て備えているようなサーバ手続きを示す。

```
server_null: sproc( void ) void = 0 ;
```

これは、何もしないサーバ手続きである。サーバが動作しているかどうかを確認するために利用される。

```
server_create_with_lower:
sproc( lowers: oid_array_t, options: option_t ) new: oid_t = 1 ;
```

これは、オブジェクトの堆積を作るサーバ手続きである。引数として下位層のオブジェクト lowers とオプション options を取り、オブジェクトを生成し、結果として生成したオブジェクトの識別子を返す。

```
server_copy:
sproc( orig:oid_t ) new:oid_t = 2 ;
```



これは、オブジェクトのコピーを行うサーバ手続きである。与えられた引数 `orig` で示されたオブジェクトに対して、いわゆる「深いコピー (deep copy)」を行う事で、新しいオブジェクトを生成する。下位層のオブジェクトがある場合、再帰的にコピーされる。結果としてその生成したオブジェクトの識別子を返す。

```
server_create_like:  
sproc( orig:oid_t ) new:oid_t = 3 ;
```

これは、`server_copy` と同様に、引数 `orig` で示されたオブジェクトと同じ形のオブジェクトを生成するサーバ手続きである。`server_copy` とは異なり、オブジェクトの内容は、コピーされない。結果として、空のオブジェクトが生成される。

```
any_server: set of sproc =  
{  
    server_null,  
    server_create_with_lower,  
    server_copy,  
    server_create_like,  
};
```

任意のサーバは、このような4つのサーバ手続きを受け付けることが要求される。

#### B-4 一般のオブジェクト

この節では、一般的なオブジェクトが全て備えているようなオブジェクト手続きを示す。これらのオブジェクト手続きは、隠された引数 `this:oid_t` で指定されたオブジェクトを操作するものである。

```
object_null: proc( void ) void = 10 ;
```

これは、何もしないオブジェクト手続きである。オブジェクトが存在するかどうかを確認するために使われる。

```
object_getattr: proc( void ) attr:attribute_t = 11 ;
```

これは、オブジェクトの属性を得るためのオブジェクト手続きである。結果として、属性を返す。属性としては、所有者、大きさ、許されたアクセス、最終更新時刻等がある。

```
object_setattr: proc( attr:attribute_t ) void = 12 ;
```

これは、属性を変更するオブジェクト手続きである。

```
object_lowers: proc( void ) lowers:oid_attr_t = 13 ;
```

これは、下位層のオブジェクトを返すオブジェクト手続きである。

```
object_kill: proc( void ) void = 14 ;
```

これは、オブジェクトを消去するオブジェクト手続きである。

```

any_object: set of proc =
{
    object_null,
    object_getattr,
    object_setattr,
    object_lowerr,
    object_kill,
};

```

任意のオブジェクトは、このような5つのオブジェクト手続きを受け付けることが要求される。

#### B-5 ファイル・サーバとファイル・オブジェクト

ここでは、ファイル・オブジェクトとしてのオブジェクト手続きとファイル・サーバとしてのサーバ手続きを示す。オブジェクト手続きの場合、隠された引数 `this:oid_t` で指定されたオブジェクトが操作対象となる。

```
file_server_null: sproc( void ) void = 50 ;
```

これは、ファイル・サーバかどうかを確認するためのサーバ手続きである。

```
file_balloc: proc( where, count: int ) void = 53 ;
```

これは、指定された領域を確保するオブジェクト手続きである。where は、確保する領域の先頭のオフセット、count は、確保する領域の長さである。これ以後、その部分に書込みが行われるまで、読み出すと0で埋められた領域が返される。

```
file_bfree: proc( where, count: int ) void = 54 ;
```

これは、指定された領域を開放するオブジェクト手続きである。where は、開放する領域の先頭のオフセット、count は、開放する領域の長さである。これ以後、この部分を読み出すと、エラーとなる。

```
file_read: proc( where, count: int ) buff: opaque = 55 ;
```

これは、指定された領域を読み込むオブジェクト手続きである。where は、読み込む領域の先頭のオフセット、count は、読み込む領域の長さである。結果として、データが格納された不定型のバッファを返す。

```
file_write: proc( where, count: int; buff: opaque ) void = 56 ;
```

これは、指定された領域に書込みを行うオブジェクト手続きである。where は、確保する領域の先頭のオフセット、count は、書込む領域の長さ、buff は、書込むデータである。

```

file_object: set of proc =
    any_object +
    {
        file_balloc,
        file_bfree,
    }

```

```

        file_read,
        file_write,
};

```

ファイル・オブジェクトは、任意のオブジェクト `any_object` に加えて、これら4つのオブジェクト手続きを受け付けることが要求される。

```

file_server_procs: set of sproc =
    any_server +
    { file_server_null } +
    file_object ;

```

ファイル・サーバは、`any_server`に指定された任意のサーバに対するサーバ手続きに加えて、サーバ手続き`file_server_null`、および、`file_object`に指定されたオブジェクト手続きを受け付けることが要求される。

```

file_server: type = server( file_server_procs );

```

`server`型構成子により、ファイル・サーバのスタブが生成される。

#### B-6 ディレクトリ・サーバとディレクトリ・オブジェクト

ここでは、ディレクトリ・オブジェクトとしてのオブジェクト手続きとディレクトリ・サーバ（名前サーバ）としてのサーバ手続きを示す。オブジェクト手続きの場合、隠された引数 `this:oid_t` で指定されたオブジェクトが操作対象となる。

```

dir_server_null: sproc( void ) void = 30 ;

```

これは、ディレクトリ・サーバかどうかを判定するためのサーバ手続きである。

```

dir_lookup: proc( name:string ) oid:oid_t = 34 ;

```

これは、ディレクトリの検索を行うオブジェクト手続きである。引数として与えられた名前 `name` をキーとして、そのオブジェクトを検索し、結果としてヒットしたオブジェクトを `oid` に返す。

```

dir_register: proc( oid:oid_t; name:string ) void = 35 ;

```

これは、そのディレクトリに、引数として与えられたオブジェクト識別子 `oid` と名前 `name` の組を登録するオブジェクト手続きである。

```

dir_remove: proc( name:string ) oid:oid_t = 36 ;

```

これは、そのディレクトリから、その名前 `name` のエントリを削除するオブジェクト手続きである。

```

dir_readdir: proc( void ) array of dir_entry_t = 37 ;

```

これは、そのディレクトリの内容を全て返すようなオブジェクト手続きである。

```

dir_object: set of proc =
  any_object +
  {
    dir_lookup,
    dir_register,
    dir_remove,
    dir_readdir,
  };

```

ディレクトリ・オブジェクトは、任意のオブジェクト any\_object に加えて、これら4つのオブジェクト手続きを受け付けることが要求される。

```

dir_server_procs: set of sproc =
  any_server +
  { dir_server_null } +
  dir_object ;

```

ディレクトリ・サーバは、any\_serverに指定された任意のサーバに対するサーバ手続きに加えて、サーバ手続き dir\_server\_null、および、dir\_objectに指定されたオブジェクト手続きを受け付けることが要求される。

```

dir_server: type = server( dir_server_procs );

```

server型構成子により、ディレクトリ・サーバのスタブが生成される。

#### B-7 メモリ・サーバとメモリ・オブジェクト

ここでは、メモリ・サーバとしてのサーバ手続きを示す。メモリ・オブジェクトとしての固有のオブジェクト手続きは、存在しない。その代わりに、メモリ・オブジェクトは、以下で示すように、ファイル・オブジェクトに対するオブジェクト手続きを受け付ける。

```

mo_server_null: sproc( void ) void = 60 ;

```

これは、メモリ・サーバかどうかを確認するためのサーバ手続きである。

```

mo_server_mo_create: sproc( size:int ) mo:oid_t = 61 ;

```

これは、メモリ・オブジェクトを生成するサーバ手続きである。この手続きは、sizeで指定された長さのメモリ・オブジェクトを生成し、結果として、生成したメモリ・オブジェクトのオブジェクト識別子を返す。

```

mo_server_getpagesize: sproc( void ) pagesize:int = 62 ;

```

これは、ページ・サイズを返すサーバ手続きである。

```

memory_object: set of proc =
  file_object ;

```

メモリ・オブジェクトは、file\_object に指定されたファイル・オブジェクトに対する手続きを受け付ける。

```

mo_server_procs: set of sproc =
  any_server +

```

```

{
    mo_server_null,
    mo_server_mo_create,
    mo_server_getpagesize,
} +
memory_object ;

```

メモリ・サーバは、any\_serverに指定された任意のサーバに対するサーバ手続きに加えて、3つのサーバ手続き mo\_server\_null, mo\_server\_mo\_create, mo\_server\_getpagesize, および、memory\_objectに指定されたオブジェクト手続きを受け付けることが要求される。

```
mo_server: type = server( mo_server_procs );
```

server型構成子により、メモリ・サーバのスタブが生成される。

## B-8 プロセス・サーバとプロセス・オブジェクト

ここでは、サーバ手続き、および、遠隔手続き呼出しにより他のプロセスを操作することが可能なオブジェクト手続きを示す。オブジェクト手続きの場合、隠された引数 this:oid\_t で指定されたオブジェクトが操作対象となる。この他に、仮想プロセッサ関係の手続きとして、同一のプロセスに対してのみ利用可能なオブジェクト手続きが存在する。それらについては、付録Cにおいて説明する。

```
process_server_null: sproc( void ) void = 70 ;
```

これは、プロセス・サーバかどうかを確認するためのサーバ手続きである。

```
process_exit: proc( status:int ) void = 71 ;
```

これは、プロセスを終了するオブジェクト手続きである。status は、終了コードである。

```
process_stop: proc( void ) void = 72 ;
```

これは、プロセスの実行を一時停止するオブジェクト手続きである。

```
process_cont: proc( void ) void = 73 ;
```

これは、一時停止したプロセスの実行を再開するオブジェクト手続きである。

```
process_getpri: proc( void ) pri:int = 74 ;
```

これは、プロセスの優先順位を得るオブジェクト手続きである。

```
process_setpri: proc( pri:int ) void = 75 ;
```

これは、プロセスの優先順位を設定するオブジェクト手続きである。

```
process_get_status: proc( void ) stat:int = 76 ;
```

これは、プロセスの終了コードを得るオブジェクト手続きである。

```
process_mo_map:
proc( mo:oid_t, addr:caddr_t; len:int; prot:prot_t ) void = 77 ;
```

これは、メモリ・オブジェクトをプロセスのアドレス空間にマップするオブジェクト手続きである。moで指定されたメモリ・オブジェクトを addr にマップする。len には、マップする長さを、prot には、保護のフラグを指定する。

```
process_mo_unmap:
proc( addr:caddr_t; len:int ) void = 77 ;
```

これは、メモリ・オブジェクトをプロセスのアドレス空間からアンマップするオブジェクト手続きである。addr から len バイトにマップされているメモリ・オブジェクトをアンマップする。

```
process_object: set of proc =
  any_object +
  {
    file_read,
    file_write,
  } +
  {
    process_exit,
    process_stop,
    process_cont,
    process_getpri,
    process_setpri,
    process_get_status,
    process_mo_map,
    process_mo_unmap,
  };
```

プロセス・オブジェクトは、任意のオブジェクト any\_object に加えて、ファイル・オブジェクトに対する file\_read, write\_write オブジェクト手続き、および、プロセス専用の8つのオブジェクト手続きを受け付ける。

```
process_server_procs: set of sproc =
  any_server +
  { process_server_null } +
  process_object ;
```

プロセス・サーバは、any\_serverに指定された任意のサーバに対するサーバ手続きに加えて、サーバ手続き process\_server\_null、および、process\_objectに指定されたオブジェクト手続きを受け付ける。

```
process_server: type = server( process_server_procs );
```

server型構成子により、プロセス・サーバのスタブが生成される。

## B-9 カーネル・サーバ群

ここでは、カーネル中に含まれるサーバを示す。

```
StdPS: process_server = 1 ;
```

カーネル中には、標準プロセス・サーバ StdPS (Standard Process Server) が存在する。これは、プロセス・サーバとしての手続きを受け付ける。

StdDS: dir\_server = 2 ;

カーネル中には、標準ディレクトリ・サーバ StdDS (Standard Directory Server) が存在する。これは、dir\_server により示されるディレクトリ・サーバとしての手続きを受け付ける。このサーバは、標準名前サーバ (StdNS, Standard Name Server) と呼ばれる。

StdFS: file\_server = 3 ;

カーネル中には、標準ファイル・サーバ StdFS (Standard File Server) が存在する。これは、file\_server により示されるプロセス・サーバとしての手続きを受け付ける。

StdMS: mo\_server = 4 ;

カーネル中には、標準メモリ・サーバ StdMS (Standard Memory Server) が存在する。これは、mo\_server により示されるメモリ・サーバとしての手続きを受け付ける。

KernelTime: server( {file\_read, file\_write} ) = 5 ;

カーネル中には、時間サーバが存在する。これは、ファイルに対する読み込み／書込みを行う手続きをインタフェースとして、現在の時刻の取得／設定を行う。

# 付録C マイクロプロセス・ライブラリの機能

ここでは、マイクロプロセス・ライブラリに含まれている主要な手続きについて述べる。

## C-1 全体の構成

マイクロプロセス・ライブラリは、第0層から第3層までの4層構造になっている。それぞれの層の主な機能は、次の通りである。

第0層では、カーネルによって提供される仮想プロセッサを実現するカーネル・コールが提供される。

第1層では、スピン・ロック (C-4. 1 参照) とコンテキスト切替え (C-4. 2 参照) を実現する機能が提供される。この層を利用することにより、各プロセッサ・アーキテクチャやバス・アーキテクチャから独立したプログラムを記述することが可能となる。

第2層では、基本的なマイクロプロセスの生成・消滅機能と、sleep/wakeup モデルに基づく簡単な同期機能が提供される。この層は、マイクロプロセス・スケジューラ (C-7 参照) とレディ・キュー・モジュール (C-8 参照) を含んでいる。これらは、応用固有のスケジューリングを行うために利用される。

第3層では、セマフォ (C-6. 1 参照)、モニタ (C-6. 2 参照)、ランデブ (C-6. 3 参照)、バイト・ストリーム (C-6. 4 参照) といった同期・通信プリミティブを提供する。これらのプリミティブのソース・コードは、各並列プログラム固有の同期・通信プリミティブを構築するプログラマにプロトタイプを与える。

## C-2 データ型

マイクロプロセス・ライブラリでは、次のようなデータ型を利用する。

env\_t

これは、マイクロプロセスのコンテキスト (環境、environment) を保持する変数ための型である。

lock\_t

これは、スピン・ロックのための型である。

struct mpcb ;

typedef struct mpcb \*mp\_t ;

これは、マイクロプロセス制御ブロック (mpcb, microprocess control block) を表す型である。マイクロプロセスを操作する時に、この型の値を指定する。mp\_t は、struct mpcb \* と同じものである。

enum mpstat

これは、マイクロプロセスの状態を表す型である。次の5つの状態がある。

M_FREE	初期
M_READY	実行可能



M_RUNNING	実行中
M_BLOCKED	待ち
M_TERMINATE	終了

usec\_t

これは、 $\mu$  秒単位の時間を表現する型である。(実際に有効な精度は、アーキテクチャによって異なる。)

### C-3 第0層 カーネル・コール

仮想プロセッサは、カーネルにより実現される。仮想プロセッサに関連したカーネル・コールには、次の4つがある。

```
int vp_allocate( n )
int n ;
```

これは、 $n - 1$  個の仮想プロセッサを新たに生成することを要求するカーネル・コールである。返り値は、仮想プロセッサの識別子である。プロセスは、生成されるとまず1個の仮想プロセッサにより実行される。そこで種々の初期設定を行った後で、このカーネル・コールを発行し並列処理を開始する。

```
int vp_sleep( t )
usec_t t ;
```

これは、現在仮想プロセッサと対応している実プロセッサを開放するカーネル・コールである。t は、開放する時間のヒントである。単位は、 $\mu$  秒である。このカーネル・コールは、応用プログラム内のマイクロプロセス・スケジューラにおいて、実行可能なマイクロプロセスが存在しなくなり、プロセッサを一時的に開放するとき用いられる。t は、単にヒントであり、正確に t  $\mu$  秒後に実プロセッサが割り当てられ、応用プロセスに実行が戻ってくる保証はない。負荷が小さい時は、指定された時間より先に戻ってくることもある。

```
vp_wakeup( n )
```

これは、vp\_switch()、あるいは、vp\_sleep() で実行を停止している仮想プロセッサの実行再開を要求するカーネル・コールである。n は、必要な仮想プロセッサの数の目安である。これは、カーネルにおいてプロセッサ割当てのヒントとして用いられる。この数より少ないプロセッサが割当てられることもあるし、逆に多くのプロセッサが割当てられることもある。ただし、このカーネル・コールを用いなかったとしても、vp\_sleep() を発行した仮想プロセッサは、指定された時間が経過した後には、実行を再開する。

```
int vp_switch( vpid )
int vpid ;
```

これは、同一プロセス内の他の仮想プロセッサへ制御を切り替えるためのカーネル・コールである。スピンロックを保持したまま実プロセッサを奪われた仮想プロセッサに制御を移すために用いられる。vpidは、次に実行すべき仮想プロセッサのヒントである。カーネルは、このカーネル・コールが発行されると、以前に実プロセッサを横取りした仮想プロセッサに対して実プロセッサを割り当てる。そのプロセスの他の仮想プロセッサも vp\_sleep()、または、vp\_switch()を発行してプロセッサを開放している場合は、自分自身で掛けたロックを待っていることを意味する。したがって、この場合はエラーである。それ以外の場合は、カーネル・コールは、同じプロセスの別の仮想プロセッサ、または、他のプロセスの仮想プロセッサにコンテキストを切り替える。

## C-4 第1層 スピンロックとコンテキスト切替え

### C-4.1 スピンロック

```
void lock_init( lp )
```

```
lock_t *lp ;
```

lp で指定されたスピン・ロック変数の初期化を行う。

```
int lock( lp )
```

```
lock_t *lp ;
```

lp で指定されたスピン・ロック変数に対して、ロックの保持を試み、成功するまで復帰しない。これは、プロセッサの test-and-set 命令や swap 命令を用いて実現されている。ただし、長い時間ロックが開放されない場合は、カーネル・コール vp\_switch() により、他の仮想プロセッサへ制御を移すことがある。

```
int lock_try( lp )
```

```
lock_t *lp ;
```

lp で指定されたスピン・ロック変数に対して、ロックの保持を試み、その結果を報告する。成功した場合、返り値として LOCK\_SUCCESS を返し、失敗した場合、LOCK\_FAILED を返す。

```
void unlock( lp )
```

```
lock_t *lp ;
```

lp で指定されたスピン・ロック変数に掛けられているロックを開放する。

### C-4.2 コンテキストの生成と切替え

```
void env_make( env, sp, pc, argv )
```

```
env_t *env ;
```

```
unsigned long int sp, pc ;
```

```
int argv[ENV_MAXARGC] ;
```

コンテキスト（環境）を生成する。spは、スタック・ポインタ、pc は、プログラム・カウンタ、argv は、軽量プロセスに与える引数である。最大 ENV\_MAXARGC 個（現在は6個）の4バイト整数が渡される。

```
int env_save( env )
```

```
env_t *env ;
```

現在のコンテキストを保存する。返り値は、0である。ただし、env\_restore() から復帰する時には、その第2引数で指定された値が返される。したがって、env\_restore() では、0以外の値を指定すべきである。

```
void env_restore( env, value )
```

```
env_t *env ;
```

```
int value ;
```

コンテキストを env で指定されたものに切り替える。これにより、制御が env\_save() で保存された所にもどる。このとき、value が env\_save() の返り値となる。この手続きは、決して復帰することがない。

## C-5 第2層 マイクロプロセスの生成と基本的な同期

## C-5.1 現在実行中のマイクロプロセス

```
private extern mp_t current ;
```

この変数は、現在実行中のマイクロプロセスを保持している。private は、仮想プロセッサごとの固有領域に割り付けられる変数であることを示している。

## C-5.2 マイクロプロセスの生成・消滅

```
mp_t mp_create( func, ssize, istat, arg0, arg1, arg2, arg3, arg4, arg5 )
```

```
void (*func)();
```

```
int ssize ; /* stack size */
```

```
enum istat ; /* M_READY or M_BLOCKED */
```

```
int arg0, arg1, arg2, arg3, arg4, arg5 ;
```

マイクロプロセスを生成する。func には、最初に実行される関数のポインタを指定する。ssize は、マイクロプロセスのスタックの大きさである。istat は、マイクロプロセスの初期状態である。M\_READY が指定され時には、mp\_wakeup() が自動的に実行される。arg0 以下は、生成されるマイクロプロセスに対する引数である。最大 ENV\_MAXARGC 個（現在は6個）まで、生成された軽量プロセスに渡される。戻り値は、mp\_t 型である。

```
mp_t mpcb_init( m, func, ssize, argv, vf )
```

```
mp_t m ;
```

```
void (*func)();
```

```
int ssize ;
```

```
int *argv ;
```

```
mpcb_vftab *vf ;
```

標準のマイクロプロセスに新たに機能を付加する時に用る、マイクロプロセスを生成する手続きである。たとえば、

```
struct my_mpcb
{
    struct mpcb    m ;
    <独自のフィールド>
};
```

のように、独自のフィールドをマイクロプロセスに付加することが可能となる。func、ssize は、mp\_create() と同じである。func には、最初に実行される関数のポインタを指定する。ssize は、マイクロプロセスのスタックの大きさである。なお、マイクロプロセスの初期状態は、M\_BLOCKED である。vf は、C++言語における virtual function の表と同様の役割を持つものである。この表には、このマイクロプロセスが終了した時に呼出される関数 mvt\_final と、デバッグ用の表示関数 mvt\_print を登録しておく必要がある。

```
struct mpcb_vftab
{
    void          (*mvt_final)();
    void          (*mvt_print)();
};
typedef struct mpcb_vftab mpcb_vftab ;
```

```
mp_exit()
```

現在実行中のマイクロプロセスを終了する。ただし、mp\_create() において func により指定された関数から復帰すると、自動的にこの関数が呼出される。

### C-5.3 マイクロプロセス間の同期

```
void mp_sleep( next, wait_for )
mp_t next ;
```

現在実行中のマイクロプロセスの実行を中断する。next は、次に実行するマイクロプロセスのヒントとして使われる。wait\_for は、どのマイクロプロセスを待っているかを表もので、実行の軌跡を得るために使われる。内部的には、current で示された現在実行中のマイクロプロセスのコンテキストを保存し、マイクロプロセスの状態を待ち状態にする。プロセッサは、マイクロプロセスのスケジューラの mp\_switch() を呼出す。これと呼出す前には、かならず現在実行中のマイクロプロセス current をロックしておく必要がある。

```
void mp_wakeup( w )
mp_t w ;
```

指定されたマイクロプロセス w を実行可能状態 M\_READY に変える。すでに実行可能状態 M\_READY あるいは、実行中状態 M\_RUNNING の時には、なにもしない。

```
void mp_yield( w )
mp_t w ;
```

マイクロプロセス w へコンテキストを切り替える。現在実行中のマイクロプロセス current は、実行可能状態になる。w に 0 が指定されると、現在実行中のマイクロプロセス以外のマイクロプロセスに制御を切替える。マイクロプロセス w は、既に実行可能状態 (MP\_READY) になっている必要がある。M\_RUNNING、M\_BLOCKED の時は、mp\_yield( 0 ) と同様、別のマイクロプロセスへ切替わる。これと呼出す前には、かならず現在実行中のマイクロプロセス current をロックしておく必要がある。

## C-6 第3層 高度な同期通信プリミティブ

### C-6.1 セマフォ

このモジュールは、計数セマフォ機能を提供する。

```
sem_t *
sem_init( sem, n )
sem_t *sem ;
int n
```

セマフォの初期化を行う。nは、計数セマフォの初期値である。

```
void sem_p( sem )
セマフォ semに対して、P命令を実行する。
```

```
void sem_v( sem )
セマフォ semに対して、V命令を実行する。
```

## C-6.2 モニタ

このモジュールは、モニタと条件変数 (condition variable) の機能を提供する。

```
mon_t *
mon_init( mon )
mon_t *mon ;
    モニタ構造体 mon の初期化を行う。
```

```
void
mon_final( mon )
mon_t *mon ;
    モニタ構造体 mon を破壊する。
```

```
void
mon_enter( mon )
mon_t *mon ;
    モニタの入口操作を実現する。
```

```
void
mon_exit( mon )
mon_t *mon ;

    モニタの出口操作を実現する。
```

```
cv_t
cv_init( cv, mon )
cv_t *cv ;
mon_t *mon ;
    条件変数 cv を初期化する。monには、対応するモニタを指定する。
```

```
cv_final( cv )
cv_t *cv ;
    条件変数を破壊する。
```

```
void cv_wait( cv )
cv_t *cv ;
    条件変数 cv に、信号が送れるまで現在実行中のマイクロプロセスを待ち状態にする。
```

```
cv_signal( cv )
cv_t *cv ;
    条件変数 cv で待っているマイクロプロセスを、実行可能状態にもどす。
```

## C-6.3 ランデブ

このモジュールは、ランデブ機能を提供する。

```
msg_t msg_create( entc )
int entc ;
    ランデブの受けポイントを作成する。entc には、受けを行うエントリの数を指定する。戻り値は、ランデブの受けポイントである。
```

```
void msg_set_receiver( msg, receiver )
msg_t msg ;
mp_t receiver ;
    ランデブの受付けポイント msg に対するランデブの受付けを行うマイクロ
    プロセスとして、receiver を登録する。
```

```
void
msg_send( msg, eno, arg_size, arg_buff, res_size, res_buff )
msg_t msg ;
int eno ;
int arg_size, res_size ;
char *arg_buff, *res_buff ;
    ランデブの呼出しを行う。eno には、呼出すエントリの番号を指定する。
    arg_size は、エントリ呼出し引数の大きさ、arg_buff は、引数の先頭番地
    である。res_size、res_buff には、それぞれエントリ呼出しの結果の大き
    さと先頭番地を指定する。
```

```
int
msg_receive( msg, arg_sizep, arb_buffp, res_sizep, res_buffp )
msg_t msg ;
    ランデブの受付けを行う。*arg_sizep、*arb_buffp には、それぞれ呼出し
    側で用意されているエントリ呼出しの引数の大きさと先頭番地が返される。
    *res_sizep、*res_buffp には、それぞれ呼出し側で用意されているエン
    トリ呼出しの結果の大きさと先頭番地が返される。関数の結果としては、エ
    ントリの番号が返される。
```

```
void
msg_reply( msg, eno )
msg_t msg ;
int eno ;
    ランデブの応答を行う。eno には、応答を行うエントリを指定する。
```

```
void
msg_enable( msg, eno )
msg_t msg ;
int eno ;
    エントリ番号 eno に対するランデブの受付けを可能にする。
```

```
void
msg_disable( msg, eno )
msg_t msg ;
int eno ;
    エントリ番号 eno に対するランデブの受付けを不可能にする。
```

#### C-6. 4 ストリーム

このモジュールは、UNIXのパイプに似たストリーム通信機能を提供する。

```
pipe_t pipe_make( size, npage )
int size ;
int npage ;
```

ストリームの通信路であるパイプを生成する。size には、パイプのためのバッファの大きさを指定する。npage には、内部的なバッファの数を指定する。

```
port_t pipe_open( pipe, mode )
pipe_t pipe ;
enum porttype mode ;    /* PORT_INPUT, PORT_OUTPUT */
    パイプ pipe をモード mode で開く。実際のデータの読み書きは、この返り値である port に対して行われる。mode には、入力モード PORT_INPUT、または、出力モード PORT_OUTPUT を指定する。
```

```
int port_read( p, buff, count )
port_t p ;
char *buff ;
int count ;
    ポート p から、buff で指定されたバッファへ count バイト読み込む。
```

```
int port_write( p, buff, count )
port_t p ;
char *buff ;
int count ;
    ポート p へ、buff で指定されたバッファから count バイト書き込む。
```

## C-7 マイクロプロセス・スケジューラ

第2層に含まれているマイクロプロセス・スケジューラは、次のような外部仕様を持っている。これらの手続きは、スケジューリングの仕組み (mechanism) を提供するものである。スケジューリングの方針は、次の節で説明するレディ・キュー・モジュールにより決定される。

```
void sched_init()
    スケジューラの初期化を行う。
```

```
void sched_main()
    main() からの呼出しポイントである。復帰することはない。
```

```
void scheduler()
    mp_exit() からの呼出しポイントである。復帰することはない。
```

```
void mp_switch( next )
mp_t next;
    mp_sleep()、mp_yield() からの呼出しポイントである。現在実行中のマイクロプロセスの実行を中断し、nextで指定されたマイクロプロセスに制御を切り替える。nextに0が指定された場合、レディ・キューより実行可能なマイクロプロセスを選択し、実行する。再びプロセッサが割り当てられた時に、この手続きから復帰する。
```

## C-8 レディ・キュー・モジュール

この節では、マイクロプロセス・スケジューラ・モジュールが使うレディ・キューモジュールの外部仕様を示す。

```
void ready_init()
```

初期化を行う。

```
void ready_enqueue( m )
```

```
mp_t m ;
```

実行可能なマイクロプロセスをキューに入れる。

```
mp_t ready_dequeue()
```

実行可能なマイクロプロセスをキューから取り出す。

```
mp_t ready_remove( m )
```

```
mp_t m ;
```

キューから指定されたmpを取り出す。

```
void ready_replace( old, new )
```

```
mp_t old, new ;
```

キューに含まれているマイクロプロセス old を new で置き換える。これは、

```
ready_remove( old )
```

```
ready_enqueue( new )
```

と等価である。

```
void ready_print()
```

デバッグ用に、キューの内容を表示する。



## 付録D プロセス操作の性能を計測するためのベンチマーク

ここでは、軽量プロセス、および、重量プロセスの性能を比較する時に用いたベンチマーク・プログラムの骨格を示す。ここであげるプログラムでは、多くの処理系で利用可能な、プロセスの生成と同期のプリミティブを利用している。

### D-1 プロセスの生成・消滅

プロセスの生成・消滅の性能を測定するために用いたプログラムの骨格を以下に示す。

```
1: run( m, n )
2:   int n, m ;
3: {
4:   int j ;
5:   for( j=0 ; j<n ; j++ ) /* outer loop */
6:   {
7:     process_create( child, m );
8:     /* create a light- or heavy- weight process */
9:     wait_for_child_exit();
10:  }
11: }
12:
13: child( rem ) /* inner loop by recursive call */
14:   int rem ;
15: {
16:   if( rem > 0 )
17:   {
18:     process_create( child, rem-1 );
19:     /* create a light- or heavy- weight process */
20:     wait_for_child_exit();
21:   }
22:   exit();
23: }
```

このプログラムは、外側の明示的なループと内側の暗黙的なループの2重ループ構造をもつ。第5行の for 文が外側のループを形成している。内側のループは、プロセス child の第18行において、自分自身を再帰的に生成することにより実現されている。この結果、内側のループでは、プロセスの入れ子が形成される。これにより、共有メモリ型マルチプロセッサにおける並列処理の効果、および、軽量プロセスやスタックのキャッシングの効果が抑えられる。逆に、外側のループの実行回数を変化させることで、キャッシングの効果を調べることができる。

このプログラムは、プロセスの生成・消滅の他に、コンテキスト切替え処理（親から子へ、子から親へ）、同期処理（子供の終了待ち）を含んでいる。

## D-2 コンテキスト切替えを伴わないプロセス間の同期

コンテキスト切替えを伴わないプロセス間の同期の性能を測定するために用いたプログラムの骨格を以下に示す。

```
1: run( count )
2:     int count ;
3: {
4:     int i ;
5:     lock_t l ;
6:     for( i=0 ; i<count ; i++ )
7:     {
8:         lock( &l );
9:         unlock( &l );
10:    }
11: }
```

このプログラムでは、1つのプロセスにおいてロックに対するロック操作／アンロック操作を行っている。このプログラムでは、ロック操作とアンロック操作を合せて1回の操作と数える。ロック／アンロックの代わりに、セマフォの場合、P命令／V命令、モニタの場合、入口操作／出口操作を用いた。

## D-3 コンテキスト切替えを伴うプロセス間の同期

コンテキスト切替えを伴うプロセス間の同期の性能を測定するために用いたプログラムの骨格を以下に示す。ここでは、セマフォ、モニタ、および、強制的にコンテキスト切替えを起こすプリミティブを用いる例を示す。ただし、共有メモリ型マルチプロセッサでは、これらのプログラムを用いたとしても、実際にはコンテキスト切替えが起こらない可能性がある。すなわち、2つのプロセスが別々の実プロセッサにより実行され、実プロセッサのレベルでは、コンテキスト切替えが起こらない可能性がある。しかしながら、同期プリミティブのオーバヘッドを測定するためには、このプログラムで十分である。

### D-3.1 セマフォを用いたプログラム

セマフォを用いたプログラムを以下に示す。

```
1: sem_t sem0 ;
2: sem_t sem1 ;
3:
4: p0( count )
5:     int count ;
6: {
7:     int i ;
8:     for( i=0 ; i<count ; i++ )
9:     {
10:        P( &sem0 );
11:        V( &sem1 );
12:    }
13: }
14:
```

```

15: p1( count )
16: {
17:     for( i=0 ; i < rcount ; i++ )
18:     {
19:         V( &sem0 );
20:         P( &sem1 );
21:     }
22: }

```

このプログラムでは、2つの私有セマフォを用いている。

### D-3. 2 モニタと条件変数を用いたプログラム

モニタと条件変数を用いたプログラムを以下に示す。

```

1:
2: monitor_t          mon ;
3: conditional_variable  cv0 ;
4: conditional_variable  cv1 ;
5:
6: p0( count )
7:     int count ;
8: {
9:     int i ;
10:    mon_enter( &mon );
11:    for( i=0 ; i<count ; i++ )
12:    {
13:        cv_wait( &cv0, &mon );
14:        cv_signal( &cv1, &mon );
15:    }
16:    mon_exit( &mon );
17: }
18:
19: p1( count )
20: {
21:    mon_enter( &mon );
22:    for( i=0 ; i < rcount ; i++ )
23:    {
24:        cv_signal( &cv0, &mon );
25:        cv_wait( &cv1, &mon );
26:    }
27:    mon_exit( &mon );
28: }

```

mon\_enter() は、モニタの入口操作、mon\_exit() は、出口操作である。cv\_wait() は、その条件が起きるまで待つプリミティブ、cv\_signal() は、条件を起こすプリミティブである。

### D-3.3 コンテキスト切替えを起こすプリミティブを用いたプログラム

コンテキスト切替えを起こすプリミティブを用いたプログラムを以下に示す。

```
1: p0( count )
2:     int count ;
3: {
4:     int i ;
5:     for( i=0 ; i<count ; i++ )
6:     {
7:         yield( p1 );
8:     }
9: }
10:
11: p1( count )
12: {
13:     for( i=0 ; i < rcount ; i++ )
14:     {
15:         yield( p1 );
16:     }
17: }
```

ここで、`yield()` がコンテキスト切替えを要求するプリミティブである。これに相当するプリミティブとしては、マイクロプロセス・ライブラリでは、`mp_yield()`、`ReSC`カーネルでは、`vp_switch()`、`SunLWP`では、`lwp_yield()`、`C-Threads`では、`pthread_yield()` があげられる。

題目： 並列分散応用プログラムを対象とした  
オペレーティング・システムの研究  
著者： 新城 靖  
日付： 1992年12月

連絡先

住所： 〒305 茨城県つくば市筑波大学  
電子・情報工学系ソフトウェア研究室  
電話： 0298-53-5163, 0298-53-6455  
F a x : 0298-53-5206  
電子メール： postmaster@softlab.is.tsukuba.ac.jp

Title: A Study on an Operating System for  
Parallel and Distributed Applications  
Author: Yasushi Shinjo  
Date: December 1992

Contact

Address: Software Lab.  
Institute of Information Sciences and Electronics  
University of Tsukuba  
Tsukuba, Ibaraki 305  
Japan  
Phone: +81 298 53 5163, +81 298 53 6455  
Fax: +81 298 53 5206  
E-Mail: postmaster@softlab.is.tsukuba.ac.jp