

オペレーティング・システムの論理記述方式の提案

田胡 和哉 益田 隆司

(筑波大学)

1. はじめに

近年、半導体技術の進歩はめざましく、それとともにハードウェア性能の向上、および、価格の低下が著しい。このため、計算機システムを、複数のプロセサにより構成し、機能、性能の向上を図ることが容易になりつつある。しかしながら、このようなシステムでは、ソフトウェアが複雑になりやすく、特に、OS (Operating Systems) は、従来の単一プロセサの場合に比べて、通信機能、スケジューリング機能等をより強化する必要があり、きわめて複雑なものになりやすい。したがって、今後、OSを容易に設計、実現する方式を開発することが、より重要な課題になるものと予想される。

OSの論理構造を明らかにすることを通じて、OS設計、実現の容易化を図ることを目的とした最初の試みは、DijkstraによるTHEシステム¹⁾である。その後、HansenのSoloシステム²⁾、HoltのTunisシステム³⁾等が開発されたが、このような設計法が広く採用されているとは必ずしも言いがたい。

これは、1) OSの論理構造が充分明確にされていないこと、2) 性能が低下しやすいこと、等の理由であると考えられる。そこで、我々は、これらの点の改善を図ることを試みている。具体的には、1) 相互排除アクセスされる資源の各々を別個のプロセサで管理し、それらを通信で結合することにより、論理的に明確な構造を持つOSを設計すること、2) 既存のOS、ここではUNIX⁴⁾とまったく同一の外部仕様を持つOSを提案方式により再設計、実現し、既存システムとの比較を行うことにより、

利点、性能の客観的な評価を行うこと、3) 評価結果をもとに改善を行い、実用性のある方式を得ること、を試みている。

今回、UNIXと同一の外部仕様を持つシステムの実現がほぼ完了したので、その設計方式、記述言語、利点、性能の概要等について報告する。

2. 設計、実現方式の概要

2.1 設計方式

提案方式によるOSは、図1. に示すよ

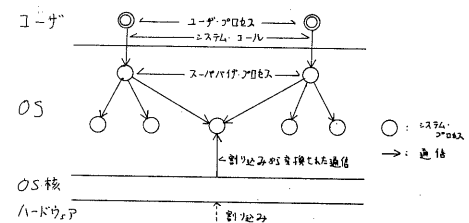


図1 提案方式によるOSの構造

うな構造を持つ。すなわち、OSをプロセス(システム・プロセス)の集合と、プロセス切り換え、割り込み処理等を実現するOS核から構成する。システム・プロセスは、ユーザ・プロセスと異なり、実行時に生成されることはなく、また、より高速に切り換えられる。各ユーザ・プロセスは、それに1対1に対応するシステム・プロセスによって管理される。これを、スーパーバイザ・プロセスと名づける。スーパーバイザ・プロセスは、ユーザ・プロセスへの資源割り当て状況から優先度を計算し、OS核に実行を依頼する。ユーザ・プロセスの発行するシステム・ゴールによる割り込みは、OS核が通信に変換し、スーパーバイザ・プロセスに伝達する。スーパーバイザ・プロセスは、その内容を分析して、要求された資源を管理す

+UNIX is a trade mark of Bell Lab. Inc.

るシステム・プロセスに対する要求に
 変換する。スーパーバイザ・プロセス以
 外のシステム・プロセスは、記憶資源、
 ファイル、ファイル入出力のためのバッ
 ファ等、相互排除アクセスされる資源
 の各々を管理しており、これらの資源
 へのアクセスは、通信によって行われ
 る。また、ディスク装置等の周辺機器も
 各々プロセスによって管理される。たと
 えば、ファイル入出力の際にOS内部で
 行われるバッファリングは、図2のよ

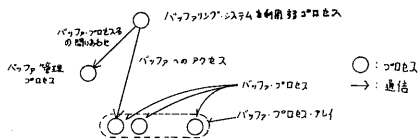


図2 プロセスによるバッファリングシステム

うにして実現される。バッファリング
 は、内部にバッファを持つバッファ
 プロセスの集合(バッファ・プロセス・
 アレイ)が実行する。バッファリング・
 システムを利用するプロセスは、バッ
 ファ管理プロセスに求める二次記憶領
 域の内容を保持しているバッファプロ
 セス名を問い合わせた後、そのプロセ
 スにアクセスする。バッファの内容は、
 通信により利用するプロセスに伝達さ
 れる。

このような設計法は、従来の、たと
 えばモニタ(monitor)を用いる方式に
 比較して、制御構造がより明確化する
 と考えられる。モニタによる方式では、
 資源の管理が手続きの系列により実現
 され、動作の中断を必要とするソフト
 ウェア割り込み等の非決定性を持つ処
 理の実現が複雑化しやすい。これに對
 して、プロセスによる方式では、資源
 管理単位に制御が独立しているのが、
 非決定的動作による影響が広範囲に広
 がりにくく、制御構造を明確化しやす
 い。

また、通信によって結合されたプロ
 セスは、変数、制御の両面において互

いに独立なもので、プロセスを単位とす
 る記述を行うことにより高いモジュラ
 リティが得られると期待される。さ
 らに、この性態は、分散処理システムへ
 の応用の際に有用である。

2.2 実現方式の概要

UNIXと同一の外部仕様を持つOSを
 実現し、従来システムと比較すること
 により、提案方式の利点、性能等を客
 観的に評価することを試みている。実
 現に際しては、比較を容易に行う立場か
 ら、従来と同一のハードウェアを用い、
 システム・プロセスはUNIXの記述言語
 であるC言語を用いて記述している。
 OS核は、アセンブリ言語で記述される。

システム・プロセス間の通信方式と
 しては、Ada言語で提案されたランデ
 ヴー(rendezvous)²⁾を用いた。これは、
 第一に、send/receive等の方式に比較
 して、同期がより明確に成立すること、
 第二に、性能向上の余地が大きいこと
 による。

ランデヴーの実行、ユーザ・プロセス
 の実行等は、システム・プロセスが手
 続き呼び出しの形でOS核に依頼するこ
 とにより実現される。この手続きを核
 プリミティブと名づける。表1に示す

通信の実行	call (呼び出し)	acc (受け付け)	end (終了)
同期制御の実行	deadlock (デッドロック判定)		
ユーザ管理	runu (ユーザの実行)	haltu (ユーザの停止)	
関連制御の実行	timeu (ユーザ・プロセスの予備時間) の管理		

表1 核プリミティブ

ように、callプリミティブによりラン
 デヴー呼び出し(rendezvous call)、
 accプリミティブにより受け付け(
 accept)、runuプリミティブによりユ
 ザの実行が行われる。

システムの記述は、プロセスを単位
 とする記述を行うと同時に、プロセス
 間の参照関係をプロセス内部の記述と
 は別個に記述する言語PNL(Process
 Network Language)を設計、実現し、
 用いている。

各プロセス内部は、C言語の手続き

の集合により構成され、通信のためのエントリ名、および、主プログラム名のみを外部に輸出する。また、ライブラリ手続き名、および、通信相手名のみを外部から輸入する。したがって、プロセス間では変数を共有せず、また、ライブラリ以外の手続きを共有しない。このため、高いモジュラリティが期待できる。さらに、このとき、プロセス内部とプロセス間の関係を分離して記述することにより、より高いモジュラリティを得ることをめざす。プロセスが、他のプロセスに対して通信を行う場合には、通信相手のエントリ名を直接輸入するのではなく、そのプロセス内でのみ使用される仮のエントリ名を定義し、その名前に対する通信として記述する。この仮の名前は、プロセス内部の記述とは独立に、PNLを用いて通信相手のプロセスのエントリに結合する。

このような記述法の利点は、第一に、プロセスを部品として扱うことにより、システムの設計、実現が容易になることである。各プロセスは、他のプロセス中で定義された名前を直接には輸入しないので、他のプロセスとは完全に独立に記述、コンパイルが可能である。このため、各プロセスの設計、実現の段階と、プロセスを組み合わせることによるシステム設計、実現の段階を分離できる。また、デバッグにおいても、各プロセスのデバッグとシステムのデバッグを分離することにより、各作業を単純、軽減化することができる。

提案方式の第二の利点は、システム構造の把握が容易になる点である。提案方式では、プロセスのみがシステムの構成要素であり、プロセスの組み合わせ構造としてシステムの構造を抽象化できると考えられる。すなわち、この組み合わせ構造を記述するものであり、ここから、システムの理解、および、

変更のために有用な情報が得られる。

これに対して、モニタによる方式では、手続き呼び出しによる抽象化、階層化を行うことにより、設計、実現の容易化を試みている。ここでは、一般に、モジュール間の関係は、外部名の輸入、輸出の陽な宣言により記述される。上位のモジュール中で、より原始的な機能を実現するモジュール中の手続き名の輸入を宣言することにより、モジュール間関係の階層性を記述している。しかしながら、並行動作システムでは、非決定性動作を実現する必要があるので、手続き呼び出しの階層性は直ちに機能の階層性に結びつかない。たとえば、端末制御プログラムは、手続き呼び出しにより、機能的により上位のユーザ・プロセス管理プログラムに対して端末からの割り込みを通知しなければならぬ。このため、外部名の輸入、輸出の宣言による方法は、非決定的動作をとともなう並行動作システムの機能の階層を記述するには不十分であり、設計、実現をむづかしくする原因になっていると考えられる。また、モニタのみの集合によるOSを実現することは不可能で、システムを構成するモジュールの性質が多様になり、かつ、モジュール間の関係が外部名の輸入、輸出の宣言、および、手続き呼び出しの形でモジュール内部にうめ込まれて記述されているので、システム全体の構造を把握することが容易でない。

提案方式は、従来の方式に比較して、以上のような利点を持つ一方、多数のプロセスが互いに関連を持って動作するためにデッドロックを生ずる危険があること、プロセス間が通信で結合されているために性能の低下をきたしやすい等の欠点がある。しかしながら、今回、これらの欠点を改善し、実用性のあるOSを実現するめどが得られたの

で、以下で報告する。

3. UNIXの再設計

2. で述べた方式にしたがひ、UNIX Ver. 6 システムを再設計した結果、図3に示すようなプロセスの配置を得た。

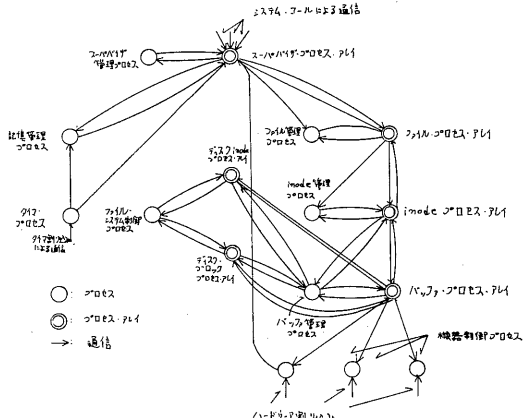


図3 再設計されたUNIXのプロセス配置

スーパーバイザ管理、記憶管理、タイマの各プロセス、および、スーパーバイザ・プロセス・アレイにより、ユーザ・プロセスへの資源割り当てをスケジュールする、いわゆるユーザ管理機能を実現している。残りのプロセスは、ファイル・システムを実現している。

スーパーバイザ管理プロセスは、スーパーバイザ・プロセスの集合（スーパーバイザ・プロセス・アレイ）に属する1つのスーパーバイザ・プロセスを選択し、ユーザ・プロセスの実行を依頼する。スーパーバイザ・プロセスは、図4(a)

に示すように記憶管理プロセスに記憶資源を要求し、その割り当てを受けた後、プログラムをロードし、CPUプリミティブを用いてユーザ・プロセスを生成する。その後、スーパーバイザ・プロセスは、システム・コールを待ち受けると同時に、図4(b)~(d)に示す方法により、ユーザ・プロセスに対する資源割り当てをスケジュールする。すなわち、(b)に示すように、端末へのアクセス等により、ユーザ・プロセスが長時間停止する場合には、記憶管理プロセスに通知する。記憶管理プロセスは、記憶資源が不足した場合、(c)に示すように、このようなユーザ・プロセス、あるいは、主記憶上に2秒以上存在しているユーザ・プロセスを管理するスーパーバイザ・プロセスに対し、スワップ・アウトの要求を行う。要求を受けたスーパーバイザ・プロセスは、スワップ・アウトを実行し、領域を解放し、再び領域が割り当てられるのを待つ。

スーパーバイザ・プロセスは、1秒およびシステム・コール毎に、haltプリミティブによりユーザ・プロセスを一時停止し、その間のプロセス使用量、入出力実行回数等から、プロセス割り当てのための優先度を計算する。

優先度計算のための時間周期は、タイマ・プロセスが通知する(d)。タイマ・プロセスは、タイマ装置からの割り込みを待ち受けた後、timeプリミティブによりユーザ・プロセスのプロセッサ使用量を計測する。また、1秒毎に、記憶管理プロセス、および、スーパーバイザ管理プロセスに時刻の通知を行う。スーパーバイザ管理プロセスは、さらに各スーパーバイザ・プロセスに通知する。

ファイル・システムは、ファイル、inode、バッファの各プロセス・アレイおよびその管理プロセス、および、二次記憶機器制御プロセス等からなる。

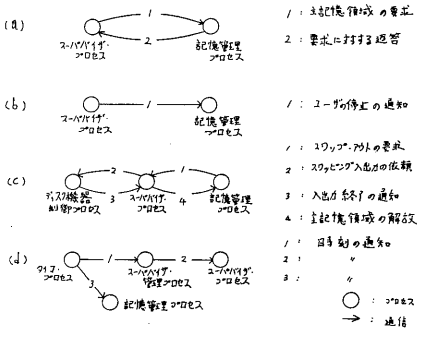


図4 ユーザ管理機能の実現

オープンされたファイルの各々を、ファイル・プロセスにより実現する。ファイル・プロセスは、ファイルのオープンに際してファイル管理プロセスが割り当てる。ファイル・プロセスは、内部に現在入出力を行ってゐるファイル部分のディスプレイメントを保持してあり、ファイル中の特定の部分の入出力をinodeプロセスに要求する。

UNIXでは、1024バイトのブロックを単位として入出力およびバッファリングが行われており、また、二次記憶上でこのブロックを単位とするマッピングを行うことにより、ファイル実体を實現してゐる。inodeプロセスは、内部にファイル実体ごとのマッピング表を保持してあり、入出力すべきディスプレイメントから二次記憶上のブロック番地を計算する。inodeプロセスは、さらに、2.で述べたバッファリングシステムに対して必要なブロックを要求し、ファイル入出力を実行する。

ディスクinodeプロセスは、二次記憶上の空きinodeを管理してあり、ファイル生成(create)に際して、inodeを1つ割り当てる。ディスク・バッファプロセスは、二次記憶上の空きブロックを管理してあり、ファイルの拡張に際して、ブロックを割り当てる。

4. 実現

現在、PDP11/34相当の計算機上で實現されたシステムが動作してあり、単一ユーザでのプログラム編集、コンパイル等が行えるようになってゐる。以下において、實現の際の問題点、その解決法等について述べる。

4.1 OS核の實現

OS核は、高速に動作することを第一の目標として設計、實現された。

Callプリミティブを、通信相手プロセス名、エントリ名、通信引数を引数として呼び出すことによりランデブ呼び出しが實現される。プロセス名は、

pcb (process control block) の番地、エントリ名は1語中のビット位置によつて識別される。accプリミティブを、受け付けるべきエントリ名を引数として呼び出すことにより、ランデブ受け付けが實現される。このとき、callプリミティブは、直ちにランデブが成立する場合にはスケジューラを介さずに直接通信を受け付けるプロセスを実行する。また、accプリミティブは、callプリミティブの引数がつまれたスタックへのポインタを値として返すことにより、きわめて高速に通信引数の結合が行われる。このため、実測の結果、ランデブ成立、および、終了後のスケジューラ等のオーバーヘッドを総計して、約150μsecでランデブが実行された。

4.2 記述言語の設計、實現

プロセス内部記述用のC言語には、entおよびcallマクロが追加された。これらのマクロ宣言から、外部名宣言が生成されるとともにcall、accプリミティブの仕様に適合する引数を生成する。

図5は、記憶管理プロセスの内部記

```

ent MFREE mfree ----- (a)
ent MALLOC malloc ----- (b)
#call SREQ sreq ----- (c)
#call ANS ans ----- (d)

process() ----- (e)
{
    int **p;
    int addr, size;

Loop:
    b = acc(MFREE : MALLOC : ---); --- (f)
    switch(b) {
        case MALLOC :
            addr = *(b + 1);
            size = *(b + 2);
            endr();
            enqueue(addr, size);
            schedule_allocation();
            goto loop; ----- (g)
        case MFREE :
            ----- (h)
            -----
    }
}

```

図5 記憶管理プロセス内部の概要

述の概要である。記憶管理プロセスは、主記憶領域の要求(MALLOC)、解放(MFREE)等の通信を受け付ける。entマクロ(a)により、外部名mfreeから、プロセス内部で用いられるエントリ名MFREEが定義されている。同様に、(c)においてスワップ・アウト要求を他のプロ

モスに対して行うための呼び出しエントリ SREQ が、外部名 steg の call マクロにより宣言されている。要求に対する返答のための呼び出しエントリ AMS も同様である (d)。

全てのシステム・プロセスは、"process" という名前を持つ手続き (e) を主手続きとしており、システムの起動時にこの手続きに制御が渡る。記憶資源の管理処理は、(f) の通信受け付けにより処理要求を受け付け、要求に対応する処理 (g) を、無限ループ (h) によりくり返すことにより実行される。

PNL は、図 6 に示すように、プロセ

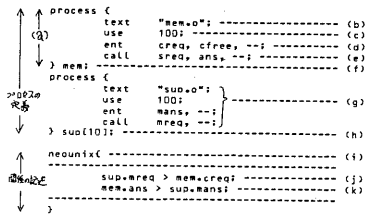


図 6 PNL 記述の例

ス定義の部分と関係記述の部分からなっている。同図では、記憶管理プロセス mem と、スーパーバイザ・プロセス・アレイ sup を定義している。(a) が、mem の定義であり、予約語 text (b) により内部記述のオブジェクト・モジュールがファイル mem.o にあることを指定される。予約語 use (c) により、プロセスが使用するスタックの量を指定する。予約語 ent (d) により、プロセス内部の ent マクロ中で宣言されている通信受け付けのための外部名を、PNL 上で宣言する。同様に、予約語 call (e) によりプロセス内部で宣言されている外部名を PNL 上で宣言する。(f) では、以上のような仕様を持つプロセスの实体を 1 個定義し、mem と名づけている。(g) では、同様にスーパーバイザ・プロセス・アレイの仕様を定義し、(h) において、同一の仕様を持つプロセスの

实体を 10 個定義し、sup と名づけている。プロセス・アレイに属するプロセスが参照する通信相手はすべて同一であるので、PNL 上では 1 個のプロセスと同等に扱う。実行時には、プロセス・アレイに属するプロセスは、アレイ名にインデックスを付けて識別する。

(j), (k) は、参照関係の記述であり、(j) において、スーパーバイザ・プロセスの記憶管理プロセスに対する領域要求、(k) において、その返答を記述している。(i) において、これらの関係により成立するシステムを、neounix と名づけている。

PNL 処理系は、以上のような PNL プログラム、および、各プロセス内部を記述した C 言語のオブジェクト・モジュールを入力として、ロード・モジュールを生成する。

4.3 システムの実現

4.3.1 ユーザ管理機能の実現

ユーザ管理機能を実現するにあたり、その問題点は、非決定性動作によるデッドロック発生の危険であった。非決定性の原因は、ユーザ・プロセス間での信号 (signal) のやりとり、端末からの割り込み、タイフ装置等に起因するものである。このような非決定性動作を、デッドロックの危険なしに実現するため、資源アクセスによりプロセスが停止する期間をなるべくみじかくするようにするとともに、常に複数の通信を待つようにした。たとえば、スーパーバイザ・プロセスがファイル・プロセスに対してアクセスする場合には、図 7 に示すように、要求と終了の通知を



図 7 ファイル入出力のための通信

別の通信により行う。ファイル入出力処理が終了する以前にも、時刻の通知あるいは、スワップ・アウトの要求等

が行われるが、スーパーバイザ・プロセスは、入出力終了の通信以外に上記の通信を待ち合わせるにより直ちに処理できる。同様の手法により、スーパーバイザ・プロセスが他のプロセスに対して通信を行う場合にも、長期間停止しないように配慮されている。

また、相互に2つのプロセスがランデブと呼び出しを行うことにより生じるデッドロックを防止するために、呼び出しを行おうとするプロセスが自分を呼び出しを待たせるか否かを判定する deadlock コリミティブを設けた。

4.3.2 ファイル・システムの実現

ファイル・システムを実現するにあたっての問題点は性能であった。ファイル入出力では、データは、バッファ → inode → ファイル → スーパーバイザ → ユーザの順に伝達される。したがって、たとえば、512 バイトのファイル入出力に対して、2048 バイトのデータ転送が必要となる。この点は、ランデブを入れ子にして用いることにより、改善を試みた。すなわち、inode プロセスは、バッファ・プロセスとのランデブ期間中にファイル・プロセスを呼び出し、同様に、ファイル・プロセスもランデブ期間中にスーパーバイザ・プロセスを呼び出すことにした。これにより、4つのシステム・プロセス間で同期が成立し、スーパーバイザ・プロセスはバッファの内容を直接ユーザ・プロセスに転送することが可能になり、データ転送の回数を1回に減らすことが可能になった。

4.4 実現結果

実現の結果、各プロセスは総計6.7k ステップ、PML コンプラムは51k ステップであった。従来のUNIXシステムは、約10k ステップであるが、より多種類の周辺機器に対応しているので、規模の点では大きな増減はない。

5. 考察

再設計の結果、第一に、モジュラリティが改善された。UNIXは、手続きの集合によって記述されており、互いに関連の深い手続きが1つのファイルにまとめられ、コンパイル単位となっている。手続き名は、他のいずれの手続きからでも参照できる。また、手続き間は、引数以外に大域変数を用いて結合されている。このような記述法によるシステムは、手続き、変数の定義、アクセス位置の把握がむづかしいために、理解、変更等が行いにくい。ここで、各コンパイル単位が輸入、輸出する外部名の数を尺度として、新旧システムのモジュラリティを比較してみる。従来のシステムのコンパイル単位であるファイルと、提案方式におけるコンパイル単位であるプロセスがそれぞれ輸入、輸出する外部名の数をシステム全体で平均すると、輸入に関して、35個から5個へ、輸出に関して9個から5個へ、それぞれ大きく減らすことができた。

従来システムとの比較は困難であるが、実現に際してのデバッグもきわめて容易に行えた。デバッグは、テスト・プロセスを用いて、機能的に下位のプロセスから順に行った。たとえば、バッファリング・システムのデバッグは、図8に示すように、バッファ・プロセス

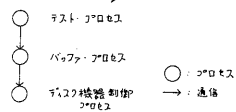


図8 バッファリングシステムのデバッグ

に通信を行い、バッファの内容を得て、それを端末に出すテスト・プロセスを作成し、バッファ、および、ディスク機器制御プロセスのデバッグを行い、以下、同様の方法で、バッファ管理、inode、inode 管理、ファイル・プロセス等の順にプロセスを追加して

スト・プロセスを作成してデバッグを行なった。テスト・プロセスは、いずれも50ステップ前後で、きわめて容易に作成することができた。この結果、ファイル・システムのデバッグは、約20人で完了した。デバッグが容易になった理由として、モジュラリティが高く各プロセスの仕様が単純なことで、各プロセスが制御的に独立していることと割り込み等の並行動作に関する処理をの核が行うので並行動作を意識する必要が少ないこと、PNLによりデバッグのための環境が容易に作成できること、等があげられる。

現段階では、性能の測定は詳細に行っていないが、表2に、新、旧システム

	UNIXでの 所要時間(sec)	新システムでの 所要時間(sec)	比(%)
J.61	147	179	73.9
J.62	240	349	67.8
J.63	161	213	73.6
J.64	279	373	74.8

表2 新旧システム上でのエンパル時間の比較

上での言語のコンパイルの所要時間を示す。この場合システムのオーバーヘッドを直接知ることはできないが、所要時間自体はる割増の増加に止まっておらず、少なくともコンパイルに関しては実用範囲内にあるものと考えられる。

デッドロックは、現在までに起きていないが、スワッピングが急増する高負荷での検査は行っていない。

実現作業の過程で見いだされた問題点として、スタック領域の不足があげられる。現在の5個のプロセス実体が存在するため、10個のスタックが必要で、スタックの使用量を厳密に評価することにより領域の節約を図っている。しかしながら、今後のプロセスは、論理空間が広がる傾向にあるので、問題は少ないものと考えられる。

6. おわりに

現在までに、提案方式の実用性に関する見通しを得た。今後の課題は、オーバーヘッドの原因を詳細に測定、解析す

ること、および、その結果必要があれば性能の改善をはかることである。

また、システムの構成要素の制御が互いに独立してゐることを生かして、分散ハードウェア用OSへの応用をはかることが、今後の課題である。

参考文献

- 1) Dijkstra, E.W. : The Structure of the "THE"-Multiprogramming System, CACM, Vol. 11, No. 5, pp. 341-346 (1968).
- 2) DoD : システム言語Ada 基準文法書, 共立出版 (1981).
- 3) Hansen, P.B. : The Solo operating system, Sofw. Pract. Exper., Vol. 6, pp. 141-205 (1976).
- 4) Holt, R.C. : Concurrent Euclid The Unix System and Tunis, Addison-Wesley (1983).
- 5) Ritchie, D.M. and Thompson, K. : The Unix Time-Sharing System, CACM, Vol. 17, No. 7, pp. 365-375 (1974).
- 6) 高平 毅 : 分散処理オペレーティング・システムの課題, 情報処理, Vol. 20, No. 4, pp. 284-288 (1979).
- 7) 田胡和哉, 益田隆司 : オペレーティング・システムの論理記述に関する考察, 情報学会第23回大会了稿集, pp. 281-282 (1981).
- 8) 田胡和哉, 益田隆司 : オペレーティング・システムの論理記述に関する考察 (第二報), 情報学会第24回大会了稿集, pp. 179-180 (1982).
- 9) 田胡和哉, 益田隆司 : オペレーティング・システムの論理構造記述法, 情報学会第25回大会了稿集, pp. 219-220 (1982).
- 10) 田胡和哉, 益田隆司 : オペレーティング・システムの論理構造記述の設計, 情報学会第26回大会了稿集, pp. 341-342 (1983).