

プロセス・ネットワークによるOSの性能について

田胡 和哉 益田 隆司
(筑波大学)

1. はじめに

近年、ハードウェア技術の発展に伴って、小型、超小型機の普及が著しく、その応用範囲も拡大している。このため、それぞれの機種、応用向きのOS、あるいは、資源を共有するための分散システム用OSがあらたに設計される機会が増加している。その一方において、OSは、設計、実現が容易でないシステムであり、その改善が種々試みられている。

このような観点から、筆者らは、OSの設計、実現の容易化を図る一つの方法として、OSを資源管理機能の集合体としてとらえ、相互排除アクセスされる資源の各々をプロセスで記述し、それらを通信を用いて結合することによりシステムを実現することを試みている¹⁾。このような通信で結合されたプロセスの集合をプロセス・ネットワークと名付けた。提案方式の狙いは、各々並行に動作するプロセスの集合によってOSを設計することにより、その制御構造を改善するとともにモジュラリティの改善を図ることにある。実際に、提案方式により、実用OSとして広く用いられているUNIX⁺システムとシステム・コール・レベルで同一の仕様を持つシステムを実現することにより、記述、デバッグ等が容易に行なえることを確認した。しかしながら、提案方式は、システムを構成する機能間の結合が通信によって実現されるために、性能が低下しやすい欠点をもつ。本稿では、実測データをもとに、従来方式との得失について解析する。

提案方式の今一つのねらいは、分散システムへの対応を図ることにある。提案方式は、その構造上、分散システムへの対応が容易である。そこで、さらに、本稿では、通信オーバーヘッドと性能の間の関係を明らかにし、提案方式による分散処理OSの性能を推定することを試みる。

2. 提案方式の概要

提案方式では、図1に示すように、OSを、通信により結合された、資源管理を行なうプロセスの集合により記述している。このとき、プロセスの切り換え、通信の実行、および、外部割込み、システム・コールの通信への変換は、OS核と呼ばれるプログラム群が実行する。ユーザ・プロセスの発行するシステム・コールは通信に変換され、OSを実現するシステム・プロセスに伝達される。このとき、システム・コールを受けつけるためにユーザ・プロセスに一对一に対応して設けるシステム・プロセスを、スーパーバイザ・プロセスと名付ける。要求の内容は、順次、より原始的な資源へのアクセスに

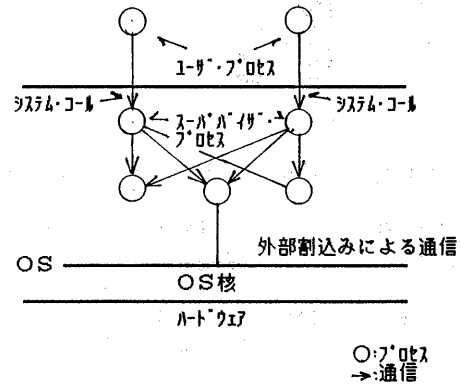


図1 提案方式によるシステムの構造

変換される。たとえば、実現したシステムでは、ファイル入出力を要求するシステム・コールは、通信により、スーパーバイザ・プロセス、オープンされたファイルの各々に割付けられるファイル・プロセス、入出力すべきディスク装置の番地を計算するinode プロセス、入出力のバッファリングを行なうバッファ・プロセス、入出力を実行するディスク・プロセスの順に伝達され、最終的に物理入出力が起動されることにより処理される。

⁺UNIX is a trademark of Bell Laboratories.

実現したシステムは、UNIXバージョン6のシステム核と同一仕様のシステム・コールを受けつけ、同一のユーザ・プログラムを実行することができる。通信方式としてランデブを用い、各プロセスの記述はC言語を用いた。それは、15種類のプロセスの、63個の実体の集合からなっており、PDP11プロセッサ上で動作する。

通信の性能はシステムの性能に与える影響が大きいので、以下で通信機構について詳しく述べる。プロセスは、手続き呼出しの形式で通信の同期に関する処理の実行をOS核に依頼する。この手続きを、核プリミティブと呼ぶ。OS核内部では、各プリミティブの呼出しに対して、状態の待避、待ち合わせ状況の判定による実行可能なプロセスの選択(スケジューリング)、その状態の回復を行なう。表1に主要な核プリ

表1 核プリミティブ

名前	機能
call	ランデブ呼び出し
acc	ランデブの受けつけ
endr	ランデブの終了

ミティブを示す。図2は、プロセス内部の通信

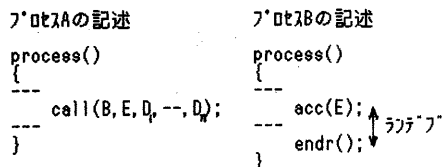


図2 通信の記述

の記述の概要である。ここでは、プロセスAがプロセスBのEエントリに対する通信を行なっている。ランデブ呼出し(rendezvous call)は、call核プリミティブを、通信相手のプロセス名(B)、エントリ名(E)、および、プロセス間で転送されるデータ(D₁~D_n)を引数として呼び出すことにより実現される。受けつけ(accept)は、accプリミティブを受けつけるべきエントリ名(E)を引数として呼び出すことにより実現される。プロセスAがcallプリミティブを

発行し、かつ、プロセスBがaccプリミティブを発行するまでどちらかのプロセスがプリミティブの副作用により停止する。両者が発行されランデブが成立すると、accは復帰してプロセスBはランデブを実行し、プロセスAはプロセスBがendrプリミティブを発行するまで停止する。プロセス間のデータ転送は、プロセスが実行する。accプリミティブは、復帰に際してcallプリミティブの引数列へのポインタを値として返す。そこで、プロセスBは、自らの局所領域へ引数を転送する。

3. 性能の予測

実現されたシステムと従来のUNIXシステムの性能の差異は、制御構造の変化に由来する原理的なものと、特に通信機構の実現方式に由来する機構的なものに分けて考えることができる。

ここで、UNIXの制御構造について考えてみる。図3は、UNIXの制御構造の概要であ

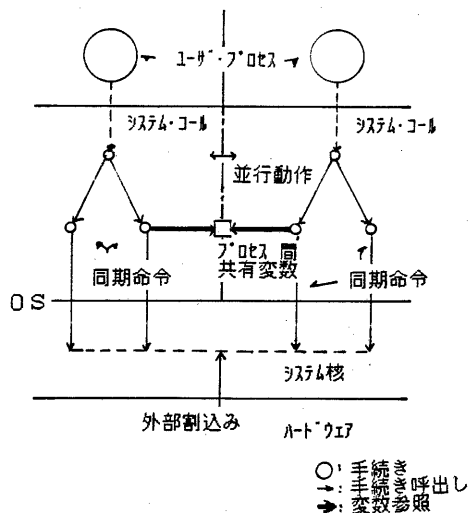


図3 UNIXの制御構造

る。システム・コールは、システム内のユーザ・プロセスに一対一に対応するプロセス中の手続き群により処理される。手続きは再入可能であり、各プロセスは同一の手続き群により動作する。ユーザ間で共有される、ファイル、主記憶領域等の資源は、プロセス間共有変数によ

て管理されており、資源アクセスの相互排除は、同期命令sleep およびwakeupを用いてプロセス間共有変数へのアクセスを相互排除することにより実現している。

ここで、制御機構の相違による性能差について考察してみる。図4に示すような、ディレクトリの探索処理の実現法について考える。この

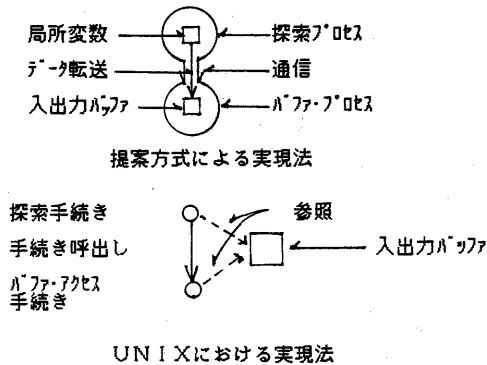


図4 ディレクトリ探索処理の実現

とき、ディレクトリは、入出力バッファ中にあるものとする。提案方式では入出力バッファはバッファ・プロセスにより実現されるので、探索を行なう探索プロセスは、通信によりバッファ・プロセスからディレクトリの内容を受け取る。これに対して、従来方式では、バッファはプロセス間共有変数により実現され、バッファ・アクセス手続きが存在する。バッファ・アクセス手続きは、他のプロセスがバッファを使用していないことを確認した後バッファ領域へのポインタを返すので、ディレクトリ探索手続きは、バッファに直接アクセスして探索を行なう。他のプロセスがバッファを使用している場合は、sleepにより待ち合わせを行なう。すなわち、通信は、バッファ・アクセス手続きの呼出し、および、sleepによる待ち合わせに対応し、機構の相違により性能が異なると考えられる。また、従来同一のプロセス中にあったディレクトリの探索とバッファの管理機能が、別個のプロセスによって実行されることによって生ずる性能の差異も考えられる。たとえば、両者は別個のループにより実現されているので、ループ制御によるオーバーヘッドが考えられる。一方バッ

ファの先読み等はより単純に実現できるので、オーバーヘッドが減少する可能性もある。

次に、今回実現した方法による通信機構のオーバーヘッドについて考察する。2.で述べたように、通信は、a)callプリミティブの引数の作成、b)状態の待避、c)実行可能なプロセスの選択、d)状態の回復、e)引数データの転送、の順に実行される。以下では、通信処理を、大きく、同期機構と(b),(c),(d)、データ転送(a),(e)、に分けて考えることにする。前者はOS核が実行し、後者はプロセスが実行する。

ディレクトリ探索の例において従来方式と比較したオーバーヘッドは、従来方式では相互排除の必要がある場合のみプロセスの切り換えが行なわれるのに対して、提案方式では通信により常にプロセスの切り換えが行なわれるので同期機構によるオーバーヘッドが生ずる。しかしながら、従来方式では、システム内のプロセスを切り換えることによりユーザ・プロセスも切り換わるので、プロセスの切り換えのたびに複雑なスケジューリングを実行する必要があり、1回の切り換えあたりのオーバーヘッドは従来方式の方が大きいと予想される。

通信によれば、従来方式よりもデータ転送量が増加する。上記の例では、従来方式においてディレクトリ探索手続きが直接バッファにアクセスしていたものが、探索プロセスの局所変数へ転送した後アクセスすることになり、あらたな転送が生じる。しかしながら、大量のデータの場合には、領域へのポインタのみ引数機構により転送し実体は局所領域間で直接転送することにより、上記a)のオーバーヘッドの減少を図っている。さらに、ファイル入出力では、入出力すべきデータはスーパーバイザ、ファイル、inode、バッファの名プロセスを経由して転送され転送量が大きくなるので、ランデブを入れ子にして用い中間のプロセスを停止してポインタのみを転送することにより、ファイル入出力における実際のデータ転送を1回に減らしている。

4. 性能の測定

実現されたシステムおよびUNIXシステム

上で同一のユーザ・プログラムを実行し、ストップ・ウォッチによるレスポンス・タイムの測定、および、ハードウェア・モニタによるOS内部の各処理の所用時間を測定した。

測定は、単一のシステム・ディスク上に実現されたOSおよびUNIXシステムを格納し、交互に動作させて行なった。このため、両システムの動作環境は、ほぼ同一である。ユーザ・プログラムとして、主に、Cコンパイラを用いた。Cコンパイラは、最も頻繁に用いられるコマンドの一つであり、典型的な負荷であると考えられる。

システムの実行特性を実時間で測定するために、ハードウェア・モニタを用いた。ハードウェア・モニタは、図5に示すように、高インピ

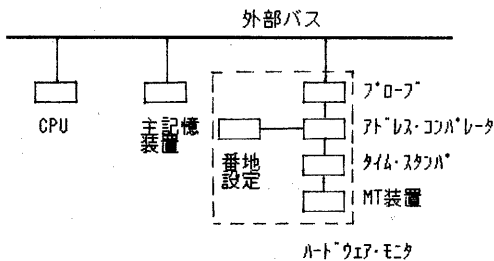


図5 性能の測定方法

ータンスのプローブによりプロセッサの外部バスの信号をモニタし、プロセッサが主記憶への書き込みを行なうたびに書き込み番地と設定番地の比較し、トリガ信号の発生、タイム・スタンプの付加、磁気テープへの記録、等を実行する。これにより、ソフトウェアが予め定めてある設定番地への書き込みを行なうと、その内容、および、時刻が磁気テープに記録されることになる。そこで、システムの状態が変化する際にこの番地への書き込みを行なう命令を付加すれば、状態の遷移を後に解析することができる。

5. 解析

5.1 性能の比較

表2は、実現されたシステムとUNIXシステム上のユーザ・プログラムの実行時間の比較であり、実行コマンドを端末に入力した後、実行が終了して次のプロンプトが印字されるまで

の時間をストップ・ウォッチで測定したものである。図中において、新システムとは実現され

表2 新旧システムの性能

ユーザ・プログラム	旧システム	新システム
Cコンパイラ ライブラリ	107 sec 26 sec	109 sec 29 sec

たシステム、旧システムとはUNIXをさし、Cコンパイラとは、Cコンパイラにより945行、12Kバイト長のソース・プログラムをコンパイルするのに要した時間であり、エディタとは、同一のソース・プログラムに対してエディタのコマンドを実行するのに要した時間であ

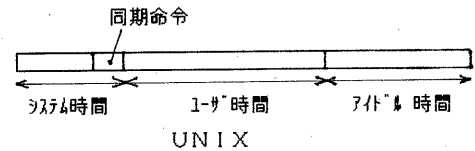
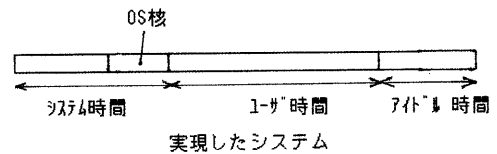


図6 処理時間の比率

る。図6は、Cコンパイラを実行した際のシステム各部のプロセッサ時間をハードウェア・モニタによって測定したものである。すなわち、システム・コールの入り口、割込み処理の入り口、アイドル手続きの入り口、および、ユーザ・プロセス・スケジューラの出口に設定番地への書き込み命令を埋め込み、ハードウェア・モニタによりシステム時間、ユーザ時間、および、アイドル時間を測定したものである。これによれば、システム時間が54%増加していることがわかる。

以下で、システム時間について考察する。予想されたように、通信が頻繁に実行されている。上記の負荷では、1112個のシステム・コールの処理に29065回の通信が実行されており、平均して1回のシステム・コールの処理に26回の通信が実行されることが測定された。また、1回にスケジュールされたシステム・プロセスの平

均実行時間は、230 μ sec にすぎず、これは、核プリミティブの処理時間にほぼ等しい。表3は、システム・プロセスが停止する要因、および、その確率を示したものであり、外部割り込みによる1.6%以外の場合には、自ら核プリミティブを発行して停止することを示している。このため、システム時間の約40%が、OS核の実行に費やされていることが判明した。これに対して、UNIXシステムでは、平均して1

表3 プロセスの停止要因

プロセスの停止要因	比率
callプリミティブの発行	24.6%
accプリミティブの発行	36.8%
endrプリミティブの発行	36.9%
外部割り込みの発生	1.6%

回のシステム・コールの処理に1.2回のプロセス切り換えが行なわれており、システム時間の約30%がsleepおよびwakeupの処理に費やされている。すなわち、提案方式によるシステムは、UNIXに比較して、同期機構では114%、その他の部分では29%それぞれプロセッサ時間が増加している。同期機構のオーバーヘッドについては、以下で解析を行なう。他の部分におけるプロセッサ時間の増加の原因は、現状では不明である。

5.2 通信機構の性能解析

はじめに、同期機構のオーバーヘッドについて解析する。表4は、OS核の各機能の、OS核

表4 OS核の各処理に要する時間

機能	実行時間	平均処理時間
callプリミティブ	26.8%	181 μ sec
accプリミティブ	47.3%	213 μ sec
endrプリミティブ	12.3%	55 μ sec
システムコール処理	2.0%	292 μ sec
外部割り込み処理	18.4%	274 μ sec

全体の実行時間に占める割合、および、1回の処理の平均時間である。これより、1回の通信の実行に要する平均時間は449 μ secであることがわかる。さらに各処理においては、状態の退避、回復にそれぞれ43 μ secかかることが判明した。そこで、通信における同期処理では、41%が状態状態の退避、回復に費やされており、

残りの59%がスケジューリングに費やされていることになる。

図7は、同期機構のオーバーヘッドとシステム性能の間の関係を明らかにするために、意図的

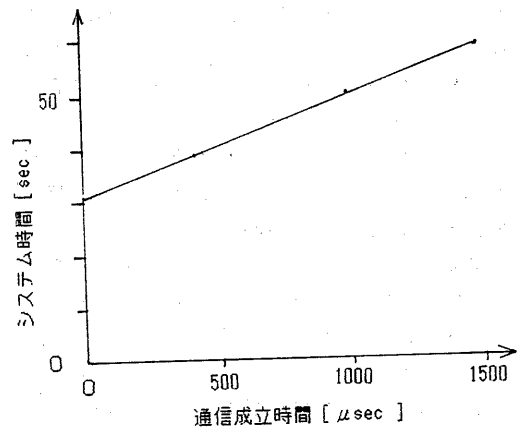


図7 同期機構の性能と全体性能

に同期機構の性能を低下させて測定を行なったものである。これは、callおよびacc核プリミティブ内部にループ命令を埋め込んで通信が成立するのに要する平均時間を延長したものであり、オーバーヘッドに比例してシステムの性能が低下しているのがわかる。

次に、通信におけるデータ転送のオーバーヘッドについて解析する。図8は、Cコンパイラが

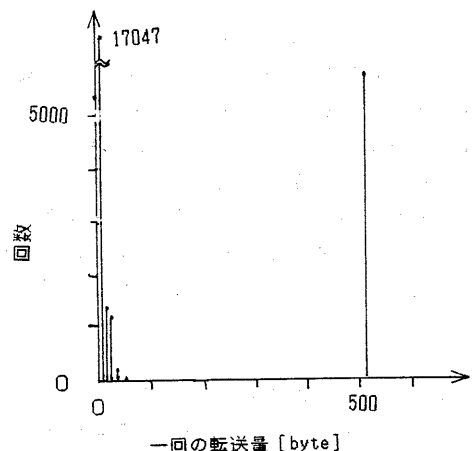


図8 データ転送量

動作している間にシステム・プロセス間で論理

的に転送されるデータ量を示したものであり、1回の通信における、callプリミティブの引数により実際の転送されたデータの転送量とポインタを用いて仮想的に転送されたデータ転送量の和を示したものである。実際の転送ののべ転送量は118Kバイト、ポインタによる仮想的なデータ転送ののべ転送量は3Mバイト、全転送ののべ転送量は3.1Mバイトであった。使用したプロセッサでは、1ワード(2バイト)の転送に5.5 μ secかかるので、実際のデータ転送に要する時間は、326 msecとなる。これは、システム時間の0.1%にすぎない。すなわち、今回の実現においては、データ転送によるオーバーヘッドの影響は少ないと言える。

しかしながら、他の通信実現方式、特に、分散システムにおいて他のプロセッサにあるプロセスとの通信では、すべて実際に転送する必要がある。そこで、同期機構の場合と同様に、ループ命令によりポインタによる転送量に比例するオーバーヘッドを付加してシステムの性能を測定した。バッファ等の連続した領域を転送する

```
1:  mov  (r0)+, (r1)+  / 転送
    sob  r2, 1b       / ループ制御
```

図9 データ転送のための機械命令

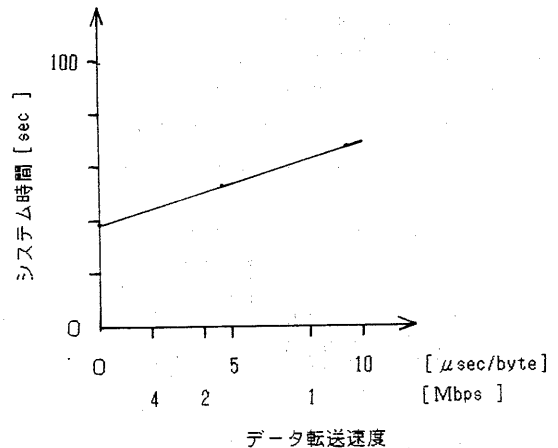
のに図9のような機械命令を用いると、1バイトあたりの転送に4.71 μ secかかる。この値に基づきオーバーヘッドを付加した測定の結果、38.4secであったシステム時間が53.3secになることが判明した。さらに、図10は、より大きな転送時間がかかると仮定して測定したシステム時間を示したものである。これにより、通信機構の性能とシステム全体の性能の関係を明らかにできる。

6. ハードウェアによる性能改善

ここでは、2台のプロセッサを交互に用いて通信における同期処理のオーバーヘッドを短縮する方法について述べる。5. で述べたように、今回の実現においては、通信機構のオーバーヘッドの大半は同期機構によって生じたものである。

そこで、ハードウェア機構によりこれを改善す

図10 転送オーバーヘッドとシステム時間の関係



ることを考える。

ここでは、2台のプロセッサA、Bを用いる。Aはプロセスを実行する。Bは、プロセスの状態の待避、実行可能なプロセスの選択、そのプロセスの状態の回復、すなわち、OS核に対応する処理を行なう。Aの実行するプロセスが通信待ち等により停止した場合には、直ちに、Bが他のプロセスを実行する。

このような方法をとることにより、同期機構のオーバーヘッドの大半を吸収することができる。また、従来のアーキテクチャによるプロセッサ、特に、ワンチップ化されたプロセッサを用いることが可能であり、容易に実現することができる。

7. むすび

通信で結合されたプロセス集合によるOSの性能について述べた。測定によれば、通信における同期機構のオーバーヘッドが大きいことが判明した。

現状では、Cコンパイラを負荷とする測定が主であるので、今後は課題は、種々の負荷に対する測定を行なって、実際の動作環境により近い測定を行なうことである。

参考文献

- 1) 田胡和哉、益田隆司：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).