

# 6V-4 分散OSの機能分割についての一検討

松下 智\*, 高畑泰志\*\*, 猪瀬 博\*, 齊藤忠夫\*  
 (\*東京大学工学部, \*\*三菱電機)

## 1. はじめに

近年, ビットマップディスプレイによるマルチウィンドウやマウスなど, 高度なマンマシンインタフェースを装備したワークステーション(WS)が注目されている. この一方では従来のTSSの持つ, ファイル・デバイス・プロセッサ・メモリなどの共有のような利点も再認識されてきた. そこで, WSにファイルの転送や共有の機能を持たせたシステムがさかんに研究されている. これらの多くは, 独立性の高いOSを持ったWSの集合体である. これに対し, 比較的新しい研究として, 通信層に異機種ノードを結合してシステム全体を一つのOSで管理し, WSの複合体上に一つの計算機を構成する方式がある. 今回提案するシステムCOSMOS/2は, これをさらに押し進めOS自体を個々に独立した論理モジュールの集合体とすることにより, 物理構成に依存せず, 自由な構成が可能な分散OSの実現への指針を得た.

## 2. システムの概要

本システムは, アプリケーションレベルでUNIXと互換性を持った分散OSをめざす. 本システムはモジュールという論理機能単位を各ノードに任意個数配置した構成をとる. ここでは, モジュールをその機能に対応して, サーバ・ハンドラ・マネージャに分類した. モジュール間の通信は, そのモジュールが同一物理ノードにある場合でも, すべてオブジェクト間通信層と呼ばれる通信層を介して統一的行なわれる.

## 3. 論理モジュールによる構成

### 3-1 物理ノードの構成

図1にディスクを持つWSノードの構成例を示す. ユーザプロセスはプロセスサーバ(PS)というモジュール内で実行される. ファイルの管理はファイルサーバ(FS)が行い, ディスク・端末・LAN装置(Ethernet)など実際の機器との入出力は, 各機器に対し論理的に一つずつ置くデバイスハンドラ(DH)が行う. 各モジュール間の通信は通信層を介してのみ行われる.

一つのノードに配置されたモジュール群は低レベルスケジューラでスケジューリングされる. これはUNIXのそれと異なり, 固定優先度を持った簡単なものであり, ノード内の割り込み処理はすべてここでされる.

こうして, ハンドラ以外のモジュールは機器依存性がなくなる. またメモリ管理は, モジュールの要求に従いメモリハンドラが行う. 通信層は, モジュール間の通信を統一につかさどり, 通信先のモジュールのロケーションをオブジェクトid(OID)より判定し, 通信のためのDHを起動するためのサブルーチンである.

ユーザプロセスは, PS内のプロセスサーバカーネル(PSK)でスケジューリングされる. ここでは, 通常のUNIXのように, プロセス寿命などを考慮した比較的高度なスケジューリングが行なわれる. PSKはまた, ユーザプロセスから出されたシステムコールの実現のため, 必要なモジュールと通信する, いわゆるプロトコルコンバータの役目を果たす. これが, おおよそ, 従来のUNIXのカーネルに相当する.

### 3-2 モジュールの種類と機能

上記, a. PS, b. FS, c. DH以外のサーバの機能を以下に示す. d~fはシステムの集中管理機能でありシステムで唯一存在する.

d. ファイルマネージャ(FM) システム全体で一意に定められた広義のファイル名を, その指し示すファイルやデバイスのid(oid)に変換する.(ネームサーバに相当)

e. プロセスマネージャ(PM) プロセスのforkの管理や, execの際, 配置すべきPSを決定したり(負荷分散の管理), シグナルの管理を行う.

f. ノードマネージャ(NM) 各ノードの動作監視やユーザ管理を行う. ここが, 動的なブートを管理する.

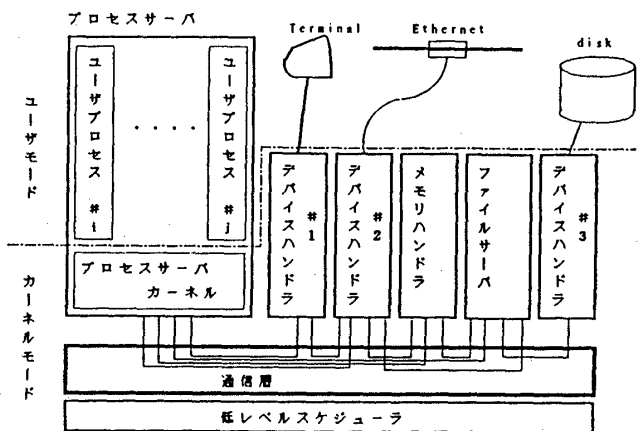


図1 ノード構成例

- g. パイプサーバ(PPS) プロセス間通信に使われるFIFO, すなわちパイプの入出力管理を行う。
- h. ローカルファイルマネージャ(LFM) ローカルディスク上でファイル名管理を行なう。
- i. ゲートウェイサーバ(GS), a~hより成るシステムを上位のネットワークなどに接続して, ツリー構成の大きなシステムにする際の接続部である。

4, オブジェクト間通信

4-1 オブジェクト識別子(oid)

本システムでは, 通信の対象を统一的にオブジェクトと呼び, システムでユニークな識別子 oid(Object ID)を与え管理する。oidは, モジュール識別子と, そこで管理される実体の識別子より構成される。本システムでは, 実体を持たないマネージャやハンドラも统一的にオブジェクトとして扱う。

4-2 通信層の構成

通信にLANを用いた場合の通信層構成を図2に示す。LLC(Logical Link Control)層では, バケットのアセンブリ, リアセンブリを行う。LLC層以下はデバイスハンドラとして実装され, 割り込みによるスケジューリングがなされる。IOC(Inter Object Communication)層は oidからどのハンドラを用いてどのノードへ通信すべきか決定をする。この構成により, LLC層以下を担当するハンドラを複数持たせ, LANを多重化することや, システムの一部を共有メモリや並列バスを用いた構成にすることが可能になる。

5 ファイルシステムの構成

本システムでは, UNIXと異なり, ディレクトリファイルのみを1つのディスクに格納し, FMで集中管理する。こうして, ファイル名を oidに変換する際, FMがいくつものFSを渡りあるく効率低下を防ぐ。この構成のファイルシステムは, ファイルアクセスが分散処理される, 1つの大きなルートファイルシステムとみなせ, UNIXと同様の方法でローカルディスクをmount, umountできる。ローカルファイルシステムはUNIXと同様の構成であり, LFM(Local FM)およびFSで管理される。こうして, カレントディレクトリがローカルにある通常のファイル処理は, ネットワークを介さずLFMのみで実行される。これを, 図3に示す。

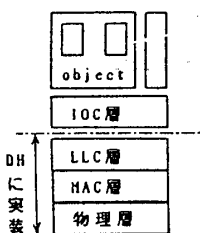


図2 通信層の論理的構成

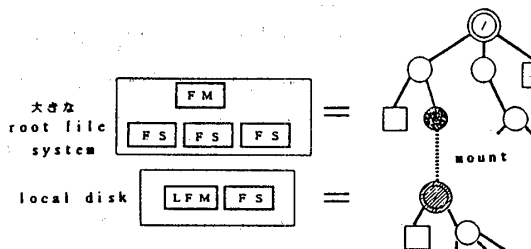


図3 ローカルファイルのマウント

6 マウントによるシステムの動的拡張

前節で述べた構成により本システムは, システムとして動作していない時にも, ローカルファイルシステムをルートとし単一ノードで立ち上がりうる。その後, 共有のFS・FMが立ち上がった時点で, ローカルなファイルシステムが自動的に共用ファイルシステムにmountされ, ファイル空間が拡張される。この点で, 独立OSを持つWSの集合体である, ホスト・ホストモデルの特徴を兼ね備えている。さらに, システムが, GSを介し上位のネットワークとつながったツリー状の構成をしていれば, 同様にしてファイル空間は動的に拡張される。この様子を図4に示す。この場合, 上層のFS, FMなどの構成は下層のものと全く同じで良く, 拡張性が高い。(プロセスに関しても同様に拡張できる。)

7. 性能予測

μVAX-II/4.2BSDのシステムコールの生起を実測し, これに従い, LANに10MbpsのEthernetを用いた構成について, 簡単な評価を行った。

この結果, μVAX程度のノード用い, それら全てをフルに動作させた場合, 87ノードでLANのボトルネックを生ずる, これは, WSを個人が占有する通常の状態では, 510台のWSからなるシステムが構成可能なことを意味する。

また, ディスク容量から一つのFMで管理できるFSの数を求めると, 87となり, 実用上充分である。

8. おわりに

資源の集中管理機能の採用により, TSSの持つ資源共有の容易さを持ち, かつ, 構成ノード数の変化に応じて共有資源の量が動的に拡張されていく分散処理システムが実現される。また, このシステムでは, OSを論理モジュールに細分化することで, 必要に応じてOS機能の分散をおこなう。さらに, これらの間の通信はオブジェクト間通信として统一的に扱い, 高い柔軟性をもつ。

現在, サーバ間で重複してアクセスするデータを極力減らし, 通信オーバーヘッドを回避するための実現法, および, ブートの実現形態などについて検討中である。

<参考文献>

[1]橋爪宏達, "ネットワーク資源共有のための統合型NOSの研究", 学位論文, 1983.

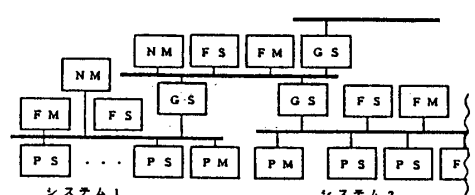


図4 階層システム構成