

マルチプロセッサにおける分散ファイル管理方式

井村 佳弘 谷口 秀夫 遠城 秀和

NTTデータ通信(株)

マルチプロセッサ用オペレーティングシステムであるDIROS (DIstributed Real-time Operating System) の分散ファイル管理方式について述べる。DIROS核の内部は機能単位でモジュール化され、各モジュール間のインタフェースにリクエスト制御と呼ぶ方式を用いている。これを利用してリモートファイルアクセス機能とリモートデバイスアクセス機能を実現し、応用プログラムのさまざまな分散形態を可能にする分散ファイルシステムを構築できることを示す。

THE DISTRIBUTED FILE MANAGEMENT METHOD ON MULTI-PROCESSORS

Yoshihiro Imura Hideo Taniguchi Hidekazu Enjo

NTT Data Communications Systems Corporation

NTT Data Yokohama Nishi Bldg., 2-11-6 Kita Saiwai,
Nishi-ku, Yokohama-shi, Kanagawa 220, Japan

The distributed file management method on multi-processors' operating system DIROS (DIstributed Real-time Operating System) is described. The kernel of DIROS is divided into some facility modules, which are connected with "request control". Remote file access and remote device access are realized by using of request control. These facilities can construct distributed file system which has various form of distribution of application programs.

1 はじめに

近年の急速なVLSI技術の進歩を背景として、複数のプロセッサを用いた価格性能比の優れた計算機や、複数の計算機をネットワーク化した分散処理システムの開発が盛んに進められている。このような分散環境を制御する分散処理OSには、各プロセッサ間での資源の共用を可能にする分散ファイル管理が必要である。このような分散ファイル管理として、これまでにRFS^[1]、NFS^[2]、DPE^[3]等が発表されている。しかし、これらの分散ファイル管理は主にファイルアクセス処理の応用プログラム（以下、APと略す）機能の分散化をねらったものであり、OS機能の分散化を行なうことはできない。そのため、サービス提供時にはその性能やメモリ効率がしばしば問題になる。

トランザクション処理用マルチプロセッサ（以降、CS：Control Stationと略す）のOSとして開発した、DIROS^[4]（Distributed Real-time Operating System）は、OS機能を分割し、そのモジュール間のインタフェースとしてリクエスト制御と呼ぶ方式を用いることにより、OS機能の分散を可能にしている。

本稿では、OS機能の分散を利用して実現したリモートファイルアクセス機能とリモートデバイスアクセス機能を持つファイル管理について、その機能と処理方式を報告する。

2 ハードウェア構成

CSのハードウェア構成例を図1に示し、その特徴を以下に述べる。

- (1) 複数のマイクロプロセッサをバスで結合している。
- (2) 共有メモリ（以降CMと略す）と割込みを利用してプロセッサ間の通信を行なう。
- (3) 各プロセッサ（以降PUと略す）の配下には、用途に合わせてディスクや通信回線

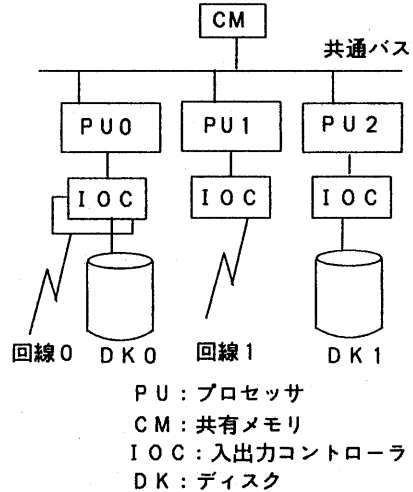


図1 ハードウェア構成例

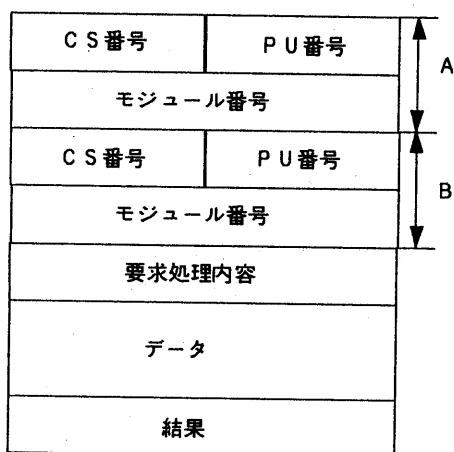
などの周辺装置を接続する。従って、ディスクを持つPUと、持たないPUが混在する。

3 DIROS構造の特徴

3.1 構造上の特徴

DIROSのプログラム構造上の特徴を以下に述べる。

- (1) 各PU上に同様なOSを分散する。
- (2) OS核の内部を機能単位でモジュール化しており、DIROS核の生成時に機能の組込みの有無を指定することができる。
- (3) OS核の各モジュール間の呼び出し（以降、内部コールと略す）を「リクエスト制御」と呼ぶインタフェースで統一している。ここで、リクエスト制御とは、各モジュール間の処理の依頼、結果応答、結果取得、取消に関する情報の受け渡しを、「リクエストブロック」と呼ぶ制御ブロックを介して行なう制御方式である。リクエストブロックの構成を図2に示す。
- (4) OSが提供するAPIインタフェース（以下、システムコールと呼ぶ）として、「非完了システムコール機能」を持つ。非完



A : 処理依頼元情報

B : 処理依頼先情報

図2 リクエストブロックの構成

了システムコール機能とは、OSへの処理依頼、結果取得、処理取消を各々別システムコールとして提供する機能である。非完了システムコール機能を用いることにより、1つのAPプロセスからOSに対して複数の処理を依頼できるため、少ないプロセス数で多くの処理を行なうことが可能となり、プロセス切り替えに伴うOSの負担を少なくすることができる。

3.2 ファイル管理の特徴

分散ファイルシステムの構築に関するDIROSファイル管理の特徴を以下に述べる。

- (1) ディスク上のファイルを通常ファイル、入出力装置や通信回線などの周辺装置を特殊ファイルと呼び、ともに木構造で管理する。
- (2) ファイルアクセスについてはキャッシュバッファ機能を持つ。APに対してキャッシュバッファを経由するアクセスと経由しないアクセスの両方を提供する。
- (3) 通常ファイルや特殊ファイルに対する排他制御機能を持つ。

4 要求条件と実現方式

4.1 要求条件

どのPU上で走行するAPプロセスからも、すべてのPUにある資源が同じように見える分散ファイルシステムを構築するためには、以下の条件が要求される。

[条件1] リモートへのアクセスとローカルに対するアクセスを同じ形式でAPに提供すること。

[条件2] ファイル管理を持たず、APプロセスが走行しないPUにある資源に対しても、他PU上のAPプロセスからアクセスできること。

4.2 実現方式

4.1節で述べた要求条件を満たす分散ファイルシステム構築のため、次の機能を実現し、全てのPUの持つ通常ファイルや特殊ファイルを木構造で一元管理できるグローバルトリーを構成する。

- (1) ファイルシステムを持つPUの通常ファイルや特殊ファイルに対するアクセスは、ファイル管理間のリクエスト制御による「リモートファイルアクセス」で実現する。
- (2) ファイル管理を持たないPUの周辺装置に対するアクセスは、ファイル管理とドライバ間のリクエスト制御による「リモートデバイスアクセス」で実現する。

リモートファイルアクセスとリモートデバイスアクセスの具体的な処理方式は5章で述べる。

図1のハードウェア構成に対し、グローバルトリーによって分散ファイルシステムを構成した場合を図3に示す。例えば、PU0上のAPから、"/PU2/f2"とファイルパス名を指定することにより、PU2上のファイル"/f2"をアクセスできる(リモートファイルアクセス)。また、PU0上の

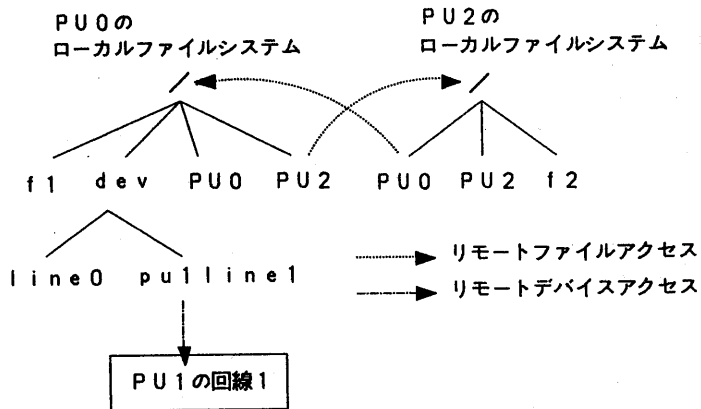


図3 グローバルトリー構成の例

APから、"/dev/pulline1"とパス名を指定することにより、PU1の回線1をアクセスできる（リモートデバイスアクセス）。

5 処理方式

5.1 リモートファイルアクセス

リモートファイルアクセスの処理方式について、ローカル/リモートの識別方式、処理の流れを以下に述べる。

5.1.1 ローカル/リモート識別方式

以下の方法により、アクセス先がローカルかリモートかを識別する。

(1) ネットワークディレクトリは、ファイルシステムが存在するCSやPUの番号に関する情報を持つ。

(2) APから指定されたパス名の中にネットワークディレクトリが含まれ、かつ、そのネットワークディレクトリのCS番号、PU番号が自PUのそれと異なる時、リモート処理と判断する。

5.1.2 処理の流れ

リモートファイルアクセス処理の例として、

PU0上のAPからPU2の通常ファイルへのリモートファイルアクセス処理の流れを図4に示し、以下に説明する。

①PU0上のファイル管理は、APから指定されたファイルパス名中のネットワークディレクトリの情報から、アクセス先がローカルかリモートかを識別する。

②PU0上のファイル管理は、処理依頼先情報としてPU2の代行プロセスを設定したリクエストブロックを作成し、その転送をリクエスト制御に依頼する。

③リクエスト制御からリクエストブロックを受取ったPU2上の代行プロセスは、その内容に基づき、依頼元のAPに代わってPU2のファイル管理に対しファイルアクセスの実行を依頼する。この時、依頼は非完了システムコール機能を利用する。

④PU2上の代行プロセスは、ファイル管理からアクセス結果を受取り、これをリクエストブロックに格納し、リクエストブロックのPU0への返送をリクエスト制御に依頼する。

⑤リクエスト制御からリクエストブロックを受けとったPU0のファイル管理は、リクエストブロック中のアクセス結果をAPに返却する。

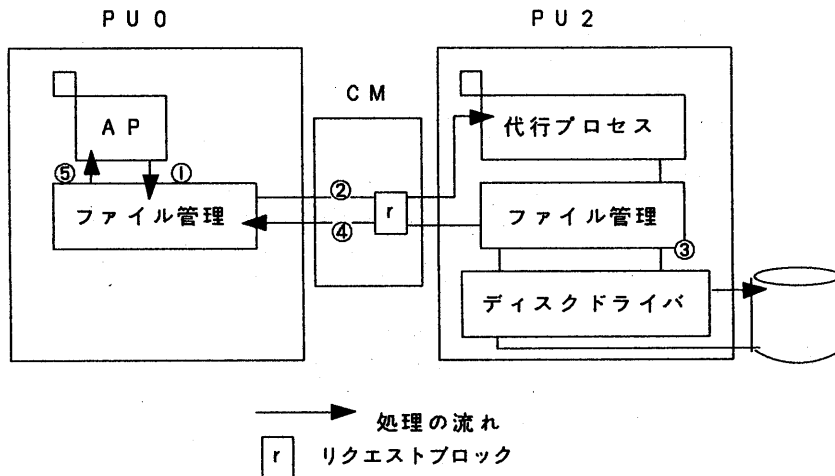


図4 リモートファイルアクセス処理の流れ

5.2 リモートデバイスアクセス

リモートデバイスアクセスの処理方式として、ローカル/リモートの識別方式、処理の流れを以下に述べる。

5.2.1 ローカル/リモート識別方式

以下の方法により、アクセス先がローカルかリモートかを識別する。

(1) 特殊ファイルは、ドライブが存在するCSやPU、およびドライブ固有の情報を「ドライブ識別子」として持つ。ドライブ識別子の内容を図5に示す。

(2) APから指定されたファイルパス名の特殊ファイル中のドライブ識別子を参照し、そのCS番号、PU番号でローカル/リモートを識別する。

5.2.2 処理の流れ

リモートデバイスアクセスの例としてPU0上のAPからPU1の回線装置へのリモートデバイスアクセス処理の流れを図6に示し、以下に説明する。

① PU0上のファイル管理は、APから指定された特殊ファイルのドライブ識別子の情報

①	②	③	④
---	---	---	---

- ① CS番号
- ② PU番号
- ③ ドライブ番号
- ④ ドライブ固有情報

図5 ドライブ識別子の内容

からアクセス先がローカルかリモートかを識別する。

② PU0上のファイル管理は、処理依頼先情報としてPU1の回線ドライブを設定したリクエストブロックを作成し、その転送をリクエスト制御に依頼する。

③ リクエスト制御からリクエストブロックを受取ったPU1上の回線ドライブは、その内容に基づき、回線装置に対してアクセス処理を実行する。

④ PU1上の回線ドライブは、アクセス処理の結果をリクエストブロックに格納し、リクエストブロックのPU0への返送をリクエ

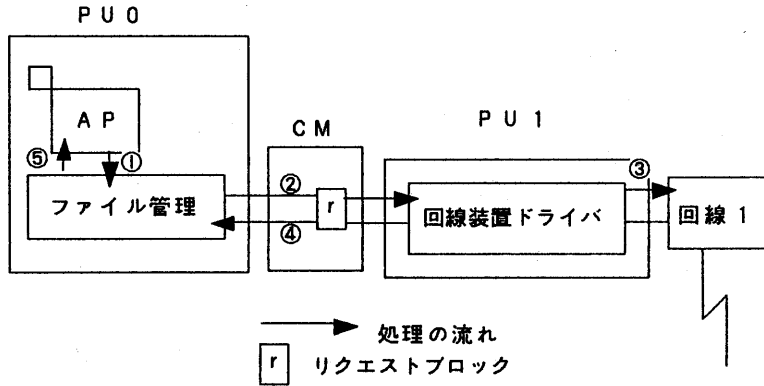


図 6 リモートデバイスアクセス処理の流れ

ト制御に依頼する。

⑤リクエスト制御からリクエストブロックを受取ったPU0のファイル管理は、リクエストブロック中のアクセス結果をAPに返却する。

6 特徴

両機能の機能上の特徴を6.1に、性能上の特徴を6.2に述べる。また、両機能による分散ファイルシステムの形態を6.3に述べる

6.1 機能上の特徴

リモートファイルアクセスとリモートデバイスアクセスの機能上の特徴比較を表1に示す。表1からわかるように、リモートデバイスアクセスは、共用、排他制御でリモートファイルアクセスに比べ機能的に劣るが、アクセス先PUのOS構成を小さくできるため必要となるメモリ量を少なくでき、ディスクの有無にかかわらず全てのPUにAPを分散できるという利点がある^[5]。

表 1 リモートファイルアクセスとリモートデバイスアクセスの機能の比較

	リモートファイルアクセス	リモートデバイスアクセス
アクセス先PUのOS構成条件	ファイル管理を必要とする	ファイル管理を必要としない
資源の共用	任意のPUで走行するAP間で資源を同時に共用することが可能	ローカルAP、リモートAP間で同じ資源を同時に共有することは不可能
資源の排他制御	任意のPUで走行するAP間で資源の排他制御が可能	ローカルAP、リモートAP間では排他制御は不可能
ディスクを持たないPU上のAPからのアクセス	不可能	ディスクを持たないPUのローカルファイルシステムを別PUのディスク上に構成することにより、可能

6.2 性能上の特徴

リモートファイルアクセスとリモートデバイスアクセスで発生するPU間通信の回数と、システムコール及び内部コールの関係を図7に示す。図7からわかるように、リモートファイルアクセスで発生するPU間通信の回数はシステムコールの発行回数に等しいが、リモートデバイスアクセスで発生するPU間通信の回数は、ファイル管理とドライバ間の内部コールの発行回数と等しい。このため、両アクセス機能の有効な適応範囲は以下のように分かれる。

(1) 通信回線などの特殊ファイルに対するアクセスでは、システムコールの発行回数と、ファイル管理とドライバ間の内部コールの発行回数が等しいため、リモートデバイスアクセスが速い。

(2) キャッシュバッファを経由した通常ファイルのアクセスでは、①キャッシュバッファへのヒット率が高いアクセス(例えば、小さなファイルを何度もアクセスするような場合)では、ファイル管理とドライバ間の内部コールの発行回数が増えるため、リモートデバイスアクセスの方が速いが、②キャッシュバッファへのヒット率が低いアクセス(例えば、大きなファイルをシーケンシャルにアクセスするような場合)では、ファイル管理とドライバ間の内部コールの発行回数が増えるため、リモートファイルアクセスの方が速い。

6.3 分散ファイルシステムの形態

リモートファイルアクセスとリモートデバイスアクセスの2つの機能を持つ分散ファイルシステムにより、性能、経済性、APの作りやすさなど、目的に合わせ、柔軟にAPやOSを構成できる。分散形態の例として、2PU間での分散ファイルシステムを性能を重視して構成した場合を図8(a)に、経済性を重視して構成した場合を図8(b)に示す。

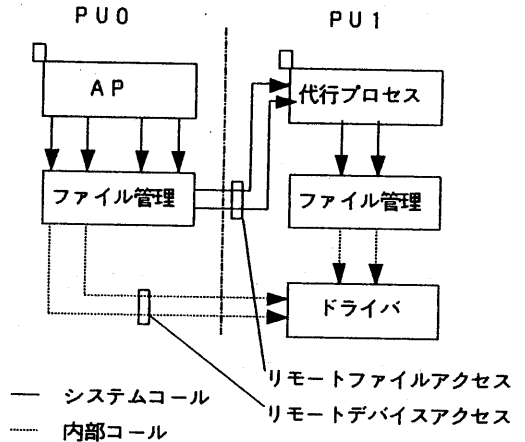
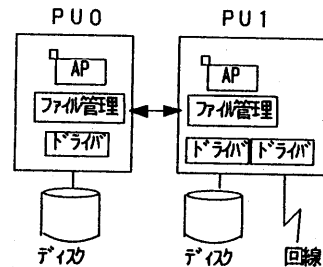
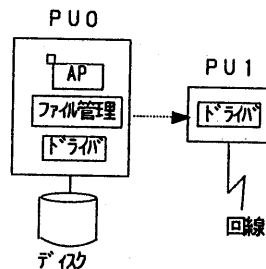


図7 PU間通信回数とシステムコール/内部コール回数の関係



(a) 性能を重視した場合



(b) 経済性を重視した場合

— リモートファイルアクセス
 - - - リモートデバイスアクセス

図8 2PU間の分散形態の例

7 あとがき

非完了システムコール機能を持ち、OS機能分散を可能にしているDIROS上で実現した分散ファイル管理方式について報告した。

ファイル管理間のリクエスト制御によるリモートファイルアクセスと、ファイル管理—

ドライバ間のリクエスト制御によるリモート
デバイスアクセスの2つの機能により、ま
まな分散形態が実現できるな分散ファイル
システムが構成できる。

〈参考文献〉

- [1] Rifkin, A. P., et al.: RFS Ar-
chitectual Overview, Proceedings
of USENIX Summer Conference,
USENIX Association, 1986
- [2] Standberg, R., et al.: Design
and Implementation of the Sun
Network Filesystem, Proceedings
of USENIX Summer Conference,
USENIX Association, 1985
- [3] 谷口、鈴木、瀬々: ファイル管理のネット
ワーク化による分散処理OSの構成法、情報
処理学会論文誌 vol. 27, No. 1, pp56-
63, 1986
- [4] 谷口: OS機能の分散を可能にするOS構
成法、電子情報通信学会論文誌D-I, vol.
J72-D-I, No. 3, pp168-174, 1989
- [5] 井村、小林、谷口: マルチプロセッサにお
けるマルチファイルシステムの構成、第38
回情報処理学会全国大会、5N-2, 1989