

分散ファイル・システムの比較と問題点

A Comparison and Problem of Distributed File Systems

中島達夫 所真理雄

Tatsuo Nakajima and Mario Tokoro

慶應義塾大学理工学部

The faculty of science and technology

Keio University

あらまし ワークステーションとローカル・エリア・ネットワークの発達により分散システムの時代へと移りつつある。分散システムを構築する際の最大の問題点はリソースの共有である。多くのOSはUnixなどのようにファイルを用いてユーザ間で情報の共有を行なっている。そのため、分散ファイル・システムの実現が分散システム構築の際の最大の問題となる。本論文では、最近の分散ファイル・システムについていくつか紹介したあと、我々が現在、開発中の分散ファイル・システムについて述べる。

ABSTRACT The age of distributed systems is coming because of the development of workstations and local area networks. In distributed systems, the most important issues is the sharing of resources. Many operating systems such as Unix make it possible to share resources through the file system. In this paper, we introduce several distributed file systems, and describe a distributed file system which we are designing and implementing.

1 はじめに

TSSの時代からワークステーションの時代への流れがますます進んでいる。以前のメインフレームと同じパワーを持つワークステーションが手頃な値段で手にはいようになり、各自が個人のワークステーションを持ちその力を最大限に利用できるパーソナル・コンピュータの実現が可能となった。そのため、そのワークステーションをネットワークでつないだ分散システムが今後の計算機の主な形態となると思われる。しかし、分散システムでは、いかにして各ワークステーションのデータを共有するかという問題を解決しなければいけない。現在、多くのOSではユーザ間のデータの共有の手段としてファイルを用いている。しかし、分散環境では各ワークステーションが個別にファイル・システムを持っているため、ファイルを共有するためにはそれらのファイル・システムを有機的につなげる分散ファイル・システムが必要不可欠である。

本論文では、現在実際に稼働中の3つの分散ファイル・システムについて述べたあと、現在、我々が開発中であるTonTos分散ファイル・システムについて述べる。

2 分散ファイル・システム

ネットワークにより複数の計算機がつけられるようになったのは15年以上前のことである。それから、実用的な分散ファイルが開発されるまで10年以上かかっている。この節では、ネットワークが開発された初期から分散ファイル・システムが登場するまでの歴史について述べる。

2.1 ファイル転送

ネットワークが開発された初期は計算機間で通信するためのプログラムとして、リモート・ログイン、ファイル転送、リモート・コマンド実行、メール・システムなどが作られた。しかし、このような環境では一般に、ユーザはリソースの存在する位置がわからないとそれにアク

セスすることが不可能だった。また、ユーザはリモート・サイトにあるファイルをアクセスするとき、ファイル転送によりいったんローカル・サイトにコピーしてからアクセスしなければならなかった。これはファイルのコピーがあちこちのサイトにできてしまうため、ファイルの一貫性と、ディスクの有効利用の面から好ましくない。

2.2 リモート・ディスク

ワークステーションに必要な高速のディスクは大きさや騒音、価格の面から机の上で用いるのには問題があるため、ディスクレス・ワークステーションが今後重要になってくると思われる。Sunではこのためにネットワーク・ディスクを導入した。これはリモート・サイトのディスクを自分のサイトのディスクのように使えるようにするためのソフトウェアである。そして、すべてのファイル・アクセスはネットワークを介してこのディスクをアクセスする。この問題点は、OSのディスク・ドライバをネットワーク・ディスクにかえただけなのでネットワークのアクセスが非常に多くなるためサーバの負荷が重くなり効率があまりよくないことである。しかし、リモート・ディスクはOSを変更することなくディスクレス・ワークステーションを実現できるので今後も普及すると思われる。

2.3 分散ファイル・システム

分散ファイル・システムは分散システム実現のキーポイントとなっている。もはや、分散ファイル・システムなしではワークステーションは存在しないとまでいえるようになってきている。初期の分散ファイル・システムではアクセスの透過性の実現が問題となった。ここでは、リモートにあるファイルもローカルにあるファイルと同じ形態でアクセスできるようにすることが最大の問題だった。その次の世代の分散ファイル・システムではファイルの位置の透過性が目標になった。つまり、ファイルがどの位置に存在していてもユーザはその位置を知らずにアクセスすることが可能となる。そして、現在の世代の分散ファイル・システムでは、ネーミングとファイル・アクセスの効率化が開発の中心となっている。次節では、この例としてNFS[1]、Andrew[2]、Sprite[3]の3つのファイルシステムについて述べる。これらは、いままでの分散ファイル・システムと異なりネーミングを工夫し、ファイル・アクセスを効率よくするためキャッシュを導入している。

3 分散ファイル・システムの現状

この節では、NFS、Andrew、Spriteの

3つの分散ファイル・システムについて、ネーミング、キャッシュの管理、ディレクトリの管理、ファイルの書き込み、トランスポート・メカニズムについて紹介し、それらの問題点について検討する。

3.1 NFS

NFSはSunが作った分散ファイル・システムである。このシステムではUnix間だけではなくMS-DOSやVMSなどヘテロジェニアスな環境でのファイル・アクセスを実現している。そのため、ファイル・アクセスのためのプロトコルを公開しそのプロトコルに従えば、どのファイル・システムのファイルでもアクセスすることが可能である。

3.1.1 ネーミング

NFSではリモート・サイトのファイルをアクセスするためリモート・サイトのサブ・ツリーをローカル・サイトにマウントする。これにより、他のサイトのファイルを自分のサイトのファイル・システムの一部としてアクセスすることが可能となる。この方式は非常に柔軟が高く様々な利用形態に対応できる。

3.1.2 キャッシュの管理

NFSではファイルの読みこみのときメモリ上にブロック・キャッシュをとることで高速化を実現している。しかし、NFSでの書き込みポリシーでは他のサイトのキャッシュの無効化まではおこなわないため、ファイルの一貫性を保つため、リードするごとにリモート・サイトからファイルを読みこまなければいけない。だが、それでは効率が悪くなるため、NFSでは一貫性を多少犠牲にして、ファイルをオープンしてからクローズするまでだけそのファイルのキャッシュを有効にしている。しかし、Unixでは一度読んだファイルのブロックを二度アクセスすることは少ないのでこの方法では効率は思ったほど向上しない。しかし、ヘテロジェニアスな環境ではキャッシュの管理を複雑にすると、プロトコルが複雑になり各OSのファイル・システムの変更が困難となるので、効率を犠牲にしてプロトコルを簡単にしている。

3.1.3 ディレクトリの管理

NFSでは、ディレクトリの解析はサーバで行なわれる。つまり、クライアントがファイルのパス名のコンポネントをサーバに渡すと、サーバでディレクトリを解析しそれが指すファイルへのポインタが返される。この方式はサーバの負荷を増やすがディレクトリのフォーマットが異なっても統一的にディレクトリを扱えるため、

異なるファイル・システムをつなぐことが可能となる。

3. 1. 4 ファイルの書出し

NFSではファイルの書出しはライト・スルーでおこなわれている。ファイルへのライトするごとにリモート・サイトがアクセスされるためサーバの負荷が非常に重くなる。また、ライト・システムコールの内部で、RPCをそのまま用いるとそれが終了するまでいつも待たされるためパフォーマンスが悪くなるので、サーバへのライトを実際におこなうのは**bi od**というデーモンにまかせ、システム・コールはすぐにリターンする。

3. 1. 5 トランスポート・メカニズム

NFSはトランスポート・メカニズムとしてRPCを用いている。このRPCでは、下位のプロトコルとしてUDPを用いている。また、エラー処理としては簡単なリトライのみしかサポートしていない。これは、NFSのサーバはステートレスなので、おなじRPCを何度繰り返してもいつも同じ結果を返すためである。また、メッセージ転送時のデータのコピーを減らすためメモリーバッファという概念を導入している。これでは、メッセージは112バイトごとのリンクド・リストとして管理され、例えばトランスポートヘッダをつけるときは、そのリンクド・リストの前にエンキューレコピーはおこなわれない。そのため、データのコピーは最小ですむ。

サーバはカーネル内のマルチ・スレッドサーバとして実現されている。そして、リクエストが来るごとにインアクティブなサーバがそれを処理する。

3. 2 Andrew

AndrewはCMUで作られた分散ファイル・システムである。このシステムはUnixの4. 2BSD上に作られ、他のOS上での使用は考えていない。このファイル・システムの特徴は数千台のワークステーション環境で用いるための分散ファイル・システムを作ることである。このシステムでは各自のワークステーションであるVirtueとファイル・サーバであるViceの2つの分類される。Vice上のファイルは単一のツリーとして管理され、そのファイルはVirtueからは同じ形態でアクセスすることが可能である。

3. 2. 1 ネーミング

Andrewでは共有ファイルはすべてVice上の単一のツリーとして管理される。そして、そのツリーは各自のワークステーションのファイル・システムにマウントすることができる。これにより、Vice上のファ

イルがローカル・ファイルのようにアクセスできるようになる。

3. 2. 2 キャッシュの管理

Andrewではファイルにアクセスするときファイル全体をローカル・ディスクへキャッシュする。つまり、ファイルをオープンしようとする、いったんファイル全体をローカル・サイトのディスクにコピーしてからアクセスされる。このようにしたのは、Unixでのファイル・アクセスはファイル全体を1度なめるようにアクセスすることが多いためである。しかし、ローカル・ディスクの存在が前提なためディスクレス・ワークステーションに対応することは不可能である。しかし、ローカル・ディスクが十分な量があれば多くのファイルがキャッシュされるのでサーバへのアクセスはほとんどおきない。しかし、この方式だと1度目のアクセスはファイル全体のコピーが必要なのでオープンの処理はファイルのサイズに比例した時間がかかる。また、ディスク容量より大きなサイズのファイルをアクセスすることはできないなどの問題が残されている。

Andrewではキャッシュの効率を向上させるためNFSより複雑なキャッシュ管理アルゴリズムを持っている。Andrewの初期のバージョンではファイルをオープンするごとにサーバにアクセスしてキャッシュが古くなっているかをチェックし、キャッシュの一貫性を保っていた。しかし、例えば、lsなどのコマンドを実行すると、何度もstatシステム・コールが実行され、そのたびごとサーバへのアクセスがおこなわれるのでパフォーマンスが悪くなる。そのため、新しいバージョンではファイルのキャッシュを持っているサイトをサーバが知っていて、どこかのサイトがファイルへ書き込みをおこなうと、そのファイルのキャッシュを持つすべてのクライアントのキャッシュを無効にする。これによりサーバの負荷が減り効率が向上した。

3. 2. 3 ディレクトリの管理

Andrewではサーバの負荷を軽くするため、ディレクトリの解析をクライアントでおこなっている。つまり、ディレクトリを読むときはクライアントにいったんディレクトリ・ファイル全体をキャッシュしてからローカル・サイトでディレクトリの解析をおこなう。しかし、すべてのディレクトリが同じフォーマットであることを前提としているため、Andrewでは異なったファイル・システムをつなぐことは不可能である。

3. 2. 4 ファイルの書き込み

Andrewではファイルの書き込みはクローズしたときに行なわれる。つまり、書き込みごとにサーバへアクセスしないので効率が向上する。しかし、コンパイル時のテンポラリ・ファイルのように、作られてすぐ消されてしまうファイルまではサーバへ書きだされるためサーバへのアクセスは思ったより減らないので効率はあまり向上しない。

3. 2. 5 トランスポート・メカニズム

Andrewではトランスポート・メカニズムとしてRPCを用いている。このRPCではNFS同様、下位のプロトコルとしてUDPを用いている。しかし、NFSと異なりサーバはステートを持っているため、複雑なエラー処理をおこなっている。

Andrewの初期のバージョンではサーバはクライアントから要求をうけとるごとに新しいプロセスを作りそれを処理していた、また、サーバ間で共有する情報はファイルに置いていた。しかし、Unixではプロセス生成のオーバーヘッド、及びファイル・アクセスのオーバーヘッドは非常に大きいのでこの方法はサーバの負荷を重くしていた。そのため、新しいバージョンでは1つのUnixプロセス上にコルーチンで動くライト・ウエイト・プロセスを実現し、要求がくるごとに新しいライト・ウエイト・プロセスを作り処理をおこなう。また、共有情報がサーバ内のメモリ上におけるのでオーバーヘッドが減り効率が向上した。しかし、Unixではコルーチンの実現は効率がよくないため、サーバはカーネル内でのマルチ・スレッド・サーバとして実現するほうが好ましい。

また、Andrewではファイル全体の転送が非常に多くなるため、ラージ・ファイルの転送の効率がよいトランスポート・メカニズムが必要となる。

3. 3 Sprite

SpriteはUCBで作られた分散ファイル・システムである。Spriteではディスクレス・ワークステーションでの使用を前提としているためキャッシュ・メカニズムが充実している。また、このシステムでは各ワークステーションは100Mバイト以上のメモリを持っていることを前提としてキャッシュ・メカニズムを設計している。この考えは、現在のメモリの価格低下を考えると好ましいと思う。この分散ファイル・システムはSprite分散OS上で動いている。

3. 3. 1 ネーミング

SpriteでのネーミングはPrefix Tableというメカニズムを用いている。これには、パス名

とそのパス名が指すファイルが存在するサーバとそのサーバ内でのオフセットが記述されている。

```
／                サーバA offset 100
／usr             サーバB offset 200
／usr/tmp        サーバC offset 150
```

パス名を解析する際、一番長くマッチするものが選び、残りのパス名をそのサーバに送る。つまり、/user/aをアクセスしようとする、これは/が一番マッチするのでサーバAにuser/aを送る。また、/usr/tmp/bをアクセスしようとする、/usr/tmp/がマッチし、サーバCにbが送られる。また、サーバが反応しないときはブロードキャストにより他のサーバにファイルのリプリケーションが存在しないかを訪ねる。そして、もし、リプリケーションが存在すると、それを持つサーバから結果が返され、それを新しいエントリとしてテーブルに登録し、それから後はその新しく登録されたサーバをアクセスする。

また、Prefix Tableではテーブルを作るのにもブロードキャストを用いている。つまり、自分が解釈できないパス名はブロードキャストを用いて他のサーバに聞く。そして、返ってきた結果をテーブルに登録する。このように、スタティックな方式にブロードキャストを組み合わせることでより柔軟な方式を実現している。

3. 3. 2 キャッシュの管理

Spriteではファイルのキャッシュはブロック単位に行なわれる。そして、キャッシュはNFSと同様メモリ上に置かれる。これは、ディスクレス・ワークステーションの使用を考えているからである。つまり、各ワークステーションが100Mバイト以上のメモリを持っていると、ユーザが使うファイルの大部分はローカル・サイトのキャッシュ上にもつてくるのが可能である。そのため、1たんキャッシュがとられるとディスク・アクセスがなくなるためファイル・アクセスの効率は非常によくなる。しかし、Spriteではファイルの先読みをおこなっていないためリードのパフォーマンスはNFSより悪い。

このシステムはキャッシュの一貫性の維持のため非常に複雑なメカニズムを持っている。今、あるクライアントがオープンしているファイルを別のクライアントがオープンしようすると、前のオープンがリード・オープンならそのファイルに関するすべてのキャッシュを無効にし、ライト・オープンならキャッシュをフラッシュしてからキャッシュ・メカニズムをストップする。この方式

では、オープンが競合するときはオーバーヘッドが大きくなるが、通常これはまれにしかおきないのでこの方式は好ましいと思われる。また、ファイルをオープンしたとき、キャッシュがローカル・サイトに存在すると、常にサーバをアクセスしてキャッシュが古くなっていないことを確認する。この方式は、Andrewの初期バージョンと同様、サーバのアクセスが多くなり効率が悪くなる。また、クライアントのみでおこなうような最適化（バス・ネーム・キャッシュなど）を用いることは不可能である。

3.3.3 ディレクトリの管理

Spriteでは自サイトでパス名が解析できなくなるとそこから先のパス名をまとめてサーバに渡す。そして、サーバがそれを解釈しその結果としてファイルへのポインタを返す。このように、Spriteではサーバのロードを低くするというをあまり考えていないのでクライアントの数が増えるとサーバの負荷が大きくなる可能性はある。

3.3.4 ファイルの書き込み

Spriteではファイルの書き込みはライト・バックで行なわれる。つまり、ファイルへのライトがおこなわれてもすぐにサーバはアクセスされない。そして、書き込みが行なわれたブロックはサーバへ30秒に1回定期的に書きだされる。これはUnixでのファイル・ブロック・キャッシュの書きこみと同じポリシーである。しかし、テンポラリー・ファイルなどはサーバに書きだす必要がないので、効率をあげるためにはこの時間を長くする必要がある。しかし、こうすると、クラッシュ時にデータが大量になくなる可能性がある。

3.3.5 トランスポート・メカニズム

Spriteではトランスポート・メカニズムとして前の2つのシステム同様RPCを用いている。しかし、このRPCでは、下位のレイヤとして他のシステムと異なりIPを用いている。また、このRPCはエラー処理を完全にサポートしている。つまり、リプライが返らないとクライアントはリクエストを再送する。そして、リクエストがいったん実行されたときはリプライは取っておかれリクエストが再送されてきたときは、そのどっぴらいたリプライを返す。

また、リクエスト、リプライのACKはできるだけビギンバックするようにして効率を向上させている。また、パラメータの性質を利用してデータのコピーを減らすようにしている。

また、サーバはNFS同様カーネル内にマルチ・スレッ

ド・サーバとして複数個作られる。これはNFSと同じメカニズムである。

3.4 分散ファイル・システムのまとめ

以上この節では3つの分散ファイル・システムについて述べた。NFSは様々なファイル・システムをつなげることができるが、効率の面で問題がある。Andrewではキャッシュをとるための時間がかかりすぎる。Spriteは分散ファイル・システムとしての完成度は非常に高いが、オープン時に必ずサーバをアクセスするのは問題がある。また、ヘテロジェニアスな環境に体してなにも考慮されていない。

4 TonTos分散ファイル・システム

TonTos分散ファイル・システムは3節で述べた分散ファイル・システムの欠点を補った分散ファイル・システムである。このファイル・システムはNFSと同時に用いることによりヘテロジェニアスな環境でも用いることが可能である。つまり、このファイル・システムを持つサイトからNFSのみ持つサイトのファイルをアクセスできるし、NFSのみ持つサイトからTonTosのファイルをアクセスすることもできる。また、NFSのサイトからのマウントはサーバに直接せず、他のクライアント・サイトにマウントし、このクライアントがサーバからキャッシュをとることによりサーバの負荷を減らすことが可能となる。

4.1 ネーミング

TonTosでのネーミングはNFSと同様、リモート・マウントにより他のサイトのファイル・システムをマウントする。この方法は、小規模なシステムから大規模なシステムまでサポートできるため非常に柔軟性が高い。しかし、NFSと異なり、サーバはサブ・ツリーを他のクライアントから容易にアクセスできるようネーム・サーバにシンボリック名として登録する。そして、クライアントは、このシンボリック名を用いてマウントする。これは、AT&TのRFSと非常に似ている。このとき、同じ名前前のサブ・ツリーが複数存在すると、それらはリアプリケーションとして扱われる。これらは、プライマリとセカンダリとに分けられ、クライアントは通常プライマリのサーバのサブ・ツリーをマウントする。しかし、このサーバがダウンしたときはセカンダリのサーバのサブ・ツリーがプライマリとなり、クライアントはこれをマウントする。これにより、アベイラビリティは向上する。クラッシュしたプライマリは立ち上がるとセカンダリとなり、新しいプライマリからファイルをコピーしてくる。

この分散ファイル・システムでは、プロセス・マイグレーションをおこなったときファイル・アクセスの効率を向上させるため、各サイトのファイル・システムのツリーの作り方を工夫している。各サイトのローカル・ファイル・ツリーは次の3つにわけられる。

- 1) Resident Tree
- 2) Nomal Tree
- 3) Temporary Tree

Resident Treeには各サイトが共通に持っているファイルが存在している。これに属するファイルはディスクレス・ワークステーションでは立上り時に全部ローカル・サイトのメモリ上にキャッシュとして持ってこられる。また、プロセスが他のサイトにマイグレーションしたとき、Resident Treeのファイルは移動した先のサイトのファイルがアクセスされる。Normal Treeは各サイトごとのローカル情報、そのワークステーションのユーザのファイル、そして、通常は使われないコマンドなどが置かれている。ディスクレス・ワークステーションでは、これらのファイルは必要に応じてサーバから持ってこられる。また、プロセスがマイグレーションしたときは、Normal Treeは移動前のサイトのものがアクセスされる。つまり、そのようなプロセスは移動前のサイトのNormal Treeへのポインタを持つことになる。Temporary Treeにはコンパイルの中間形などが置かれる。つまり、このツリーにあるファイルはキャッシュがあふれないかぎりサーバに書きだされない。ここにあるファイルは、メモリ上に作られ、ディスク・アクセスされずに消される。

4.2 キャッシュの管理

TonTosではSpriteと同様にブロック単位のキャッシュをメモリ上におこなう。このキャッシュはそのファイルのID（ネットワークID+ローカル・ファイルID）とそのファイル内のブロック番号からなる。このブロックはLRUアルゴリズムに基づきおいだされる。また、間接参照を持つブロックのキャッシュも重要である。このファイル・システムも各ワークステーションが100Mバイト以上のメモリをもっていることを前提としている。また、ディスクレス・ワークステーションを効率よくサポートするためResident Treeに関してはシステムのブート時にバルク・データ転送プロトコルを用いてローカル・サイトのキャッシュに転送される。また、ディレクトリの解析はAndrew同様クライアントでおこなわれる。

TonTosでのキャッシュの一貫性の維持は、ファイルへのライトはリードより少ないのを利用して、ライトよりリードのパフォーマンスを重視する。TonTo

sではリード時はサーバにアクセスしない。しかし、ライト・オープンの際は必ずサーバへそれを知らせる。サーバはそれにより、キャッシュを持っているサイトを知っているのこのときはすべてのサイトのキャッシュをストップする。また、このとき、キャッシュの内容が変更されていたらそれをフラッシュしてからキャッシュをストップする。この方法はリード時にはサーバへのアクセスが必要ないので効率向上する。

4.3 ファイルの書きこみ

TonTosではファイルの書きこみ戦略としてライト・バックを用いている。そして、Spriteと同様に定期的にダーティ・ブロックの書きだしをおこなう。これはNormal Treeに対してのみおこなう。Resident Treeの変更に関しては一貫性の管理など問題が残されているがこれについては今後の課題である。しかし、現在作成中のプログラムなどはエディティングとコンパイル、実行のくりかえしなので、このようなファイルは定期的な書きこむ時間を長くすれば効率がよくなる。しかし、クラッシュしたときは情報が大量になくなる可能性があるためこのへんのトレードオフが難しい。

4.5 トランスポート・メカニズム

TonTosではトランスポート・メカニズムとしてRPCを用いている。このRPCはこのファイル・システム専用のもので、このために最適化している。つまり、ファイルのアクセスはオープンしてから、クローズするまで複数回のリード・ライトがおこなわれる。そのため、パーチャル・サーキットをベースにしたRPCが好ましい。つまり、このRPCでは最初のリクエストによりコネクションが設定される。そのため、リクエスト、リプライのアックは次のリクエスト、リプライのメッセージにビギンバックすることが可能となる。このコネクションはサーバが一定時間リクエストがこないのを検出してコネクションを解放する。そして、次のリクエスト時に再び新たなコネクションを設定する。このようにすることにより、コネクションの設定時にパラメータをネゴシエートすることによりフローコントロールなどをおこなうことができネットワーク、及びサーバの負荷を軽くすることができる。

また、ワークステーションの立上り時にResident Treeをすべて転送するためバルク・データ転送プロトコルが必要である。これは、大量データ転送専用の効率のよいプロトコルである。

TonTosでのサーバはNFS、Sprite同様カーネル内のマルチ・スレッド・サーバとして作られる。

5 分散ファイル・システムから分散OSへ

TonTosはUnix上の分散ファイル・システムではなくTonTos分散OS上の1コンポーネントとして最終的には作る予定である。この節ではTonTos分散OSの他の特徴について述べる。

5.1 仮想記憶

Unixでは仮想記憶のスワップ・エリアはローカル・ディスク上にファイル・システムを介さず直接アクセスしている。しかし、TonTosでは仮想記憶はファイル・システムを介してアクセスする。つまり、各プロセスのスワップ領域はファイルとして作られる。そのため、スワップエリアを容易にリモート・サイト上に作ることができるため、プロセス・マイグレーションの実現が容易となる。

5.2 プロセス管理

TonTosでのプロセスはUnixのプロセスと違いライトウエイト・プロセスをサポートしている。これはスタック・エリア以外の領域、及びオープン・ファイル情報などの各プロセスの環境もすべて共有しているプロセスである。また、プロセスがフォークするときの効率を向上させるためCopy-on-writeを用いている。

TonTosのプロセス管理の最大の特徴はプロセス・マイグレーションのサポートにある。プロセス・マイグレーションで重要なのはプロセスが移動してもファイル・アクセスの環境が変わらないことである。プロセスがマイグレートされても、スワップエリアはファイルとして管理されるため分散ファイルを介して簡単にアクセスできる。このとき、マイグレーションを行なう前にすべてのダーティ・ページをスワップ・ファイルに書きだし、次にプロセスの環境をすべて移動先のサイトの転送する。そのとき、オープン中のファイルはリモート・サイトへそのファイルのIDや現在のオフセットを転送する。オープン中のテンポラリ・ファイルはマイグレーションがおきるとサーバに書き出され、新しいサイトはそれをサーバから読み出す。また、新たにオープンするファイルはNormal Treeに関しては移動前のサイトのものをアクセスする。Resident Tree、Temporary Treeに関しては移動先のものが用いられる。

5.3 IPC

分散OSではネットワーク透過のIPCが非常に重要である。TonTosではネットワーク透過なポートを用いてこれを実現している。これはMach [4] で用いられているものほとんど同じである。これを用いることにより各プロセスは相手プロセスが存在するサイトを知らずに通信することが可能となる。

6 結論

本論文では現在、稼動中の3つの分散ファイル・システムについて比較、検討をおこなった。そして、これらの欠点を改良したTonTos分散ファイル・システムについて提案した。これは我々が現在設計中の分散OS: TonTosの1コンポーネントである。そして、最後はTonTos分散OSについて簡単に述べた。

- [1] R.Sandberg et al. "Design and Implementation of Sun Network File System" Proc. Summer USENIX 1985
- [2] J.H.Morris et al. "Andrew: A Distributed Personal Computing Environment" CACM Vol 29 No 3
- [3] M.Nelson et al. "Caching in the Sprite Network Filesystem" Proc. SOSP 1987
- [4] M.Accetta et al. "A New Kernel Fundaton for Unix Development" Proc. Summer USENIX 1986