

V 分散システム

下島 健彦
日本電気(株)

スタンフォード大学で開発されている V 分散システムについて紹介する。V システムは、分散カーネルとサーバの集合により構成される。分散カーネルがネットワーク透過な通信を提供し、それ以外の OS 機能を通信によって起動されるサーバで実現することで、ネットワーク透過な OS を実現している。ネットワーク・プロトコルとして、VMTP と呼ばれる RPC を効率よくサポートするプロトコルを採用していること、通信の速度を重視した数々の最適化を行なっていること、プロセスのグループに対する通信(マルチキャスト)を提供していること、分散した名前管理を行なっていることなどが特徴である。

The V Distributed System

Takehiko Shimojima
NEC Corporation

This paper is an introduction of the V distributed operating system. The V consists of a distributed kernel and user level servers. The kernel provides network transparent interprocess communication (IPC) facility, process and memory management and device drivers, while servers provide other operating system functions invoked through IPC. The V uses the VMTP transport protocol which is optimized for request-response behavior, makes several optimizations for fast IPC, provides multicast to a group of processes, uses decentralized naming management.

1 概要

V分散システムは Cheriton 教授の指導のもとに1980年代初めからスタンフォード大学で開発されている分散システムで、現在も拡張や新しいアイデアの実験環境として開発が続けられている。著者は1988年夏から1年間スタンフォード大学に留学する機会を得、Vシステムプロジェクトに参加した。本論文ではその経験を元に、Vシステムの紹介を行なう。

1.1 設計思想

文献 [3] に述べられているように、Vシステムでは次の三つの設計思想が掲げられている。

1. 高速な通信の重要性

分散システムは高速な通信を基礎として、(そのコスト意識の上に)設計されるべきである。また、通信はスループットだけでなく、レスポンス性能を重視したトランザクション・オリエンテッドなものが必要である。

2. プロトコルの重要性

ネットワーク・プロトコルだけでなく、I/O、ネーミング、認証などのプロトコルがシステムを規定する。大学での研究の目的は、そういったプロトコルを開発することであり、システム作りは、開発したプロトコルを検証するためのものと考えている。

3. 小さい分散カーネルとサーバによるシステム構成

通常のOS機能は、比較的小さな分散カーネルとプロセスレベルのサーバにより実現できる。分散カーネルはネットワーク透過な通信とプロセスの走行環境であるプロセス管理、メモリ管理を実現する。

VシステムではOSの機能をサーバで実現し、サーバへの処理要求をネットワーク透過な通信によって行なうことで、OS機能をネットワーク透過にしている。

1.2 アプリケーション分野とハードウェア環境

分散カーネルは特定のアプリケーション分野に依存せず、汎用に設計されている。例えばプロセス・スケジューリングもプリエンティブ固定プライオリティの単純なものが提供されており、リアルタイムシステムもTSSのような環境もサポートできるようにになっている。

サーバやコマンドを含んだVシステム全体としては、主にプログラム開発などを行なうワークステーションのようなアプリケーションが想定されており、エディタ、パラレルmake、デバッガなどのコマンド群が用意されている。さらに、Vシステムをファクトリ・オートメーション(FA)に応用した例もあるとのことである。

ハードウェア環境としては現在、Sun2, Sun3, μ VAX-II, DEC3100(RISCワークステーション), DEC Firefly(密結合型マルチプロセッサ), VMP [5] (Cheriton教授のもとで開発されているソフトウェア制御のキャッシュを持った密結合型マルチプロセッサ)等が一つのネットワーク(イーサネット)に混在できるヘテロジニアスな環境がサポートされている。ユーザインタフェースとしてはビットマップとマウスが想定されている。

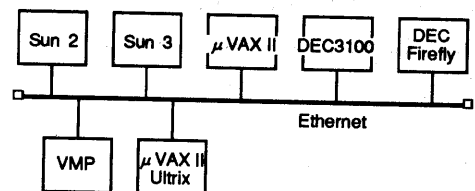


図1: ハードウェア構成

2 分散カーネル

分散カーネルはネットワーク透過なプロセス間通信 (IPC) とプロセス管理, メモリ管理, デバイスドライバなどを実現している。このうち, IPC 以外は全てサーバ (カーネルサーバと呼ばれる) の形で実現されている。カーネルサーバについては 2.7 で述べる。

2.1 V の IPC

V システムでは図 2 のような同期型の IPC が提供されている。V の IPC の特徴として次の点があげられる。

- RPC に適した単純なプリミティブ。
- 固定長メッセージ (32 バイト) + 可変長データ・セグメント。
- RPC に適した通信プロトコル VMTP。
- 通信のためのカーネルの最適化。
- マルチキャストのサポート。

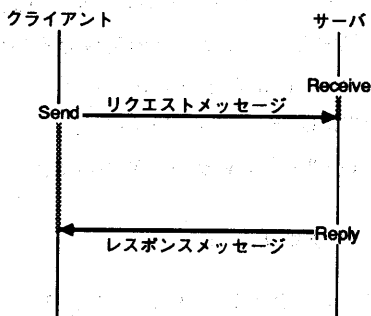


図 2: V システムの IPC

2.2 IPC プリミティブ

IPC プリミティブは Send-Receive-Reply を基本としている。特に Send は要求の送信とレスポンスの受信という RPC での通常ケースを 1 プリミティブで行なうため, オーバーヘッドが少い。また, クライアントが指定した送信バッファにレスポンスを直接書き戻すことでバッファ管理も単純化している。

2.3 固定長メッセージ

IPC は 32 バイトの固定長メッセージを基本とすることでバッファ管理を単純化している。32 バイト以上のデータの送受信が必要な場合は可変長のデータセグメントを付加する。実際に V システムでメッセージ・トラフィックパターンを調査した結果 [9], 5 割が 32 バイト以下のデータであることが報告されている。

2.4 VMTP

VMTP [4, 10] は V プロジェクトで開発されたトランスポート層のプロトコルで次のような特徴を持つ。

- コネクションの設定 / 切断がなく, リクエスト / レスポンス主体の通信に向いている。
- レスポンスがリクエストの Ack を兼ね, リクエストが直前のレスポンスの Ack を兼ねることで最小 2 パケットで RPC を実現。
- ヘッダ中に固定長メッセージが入る。
- マルチキャスト, フォワーディング, ストリーム, プライオリティなどの豊富な機能。
- 選択的再送, レートベースのフロー制御など, 高速, 低エラー率のネットワークへの対応。
- Idempotent なレスポンスの最適化。

再実行しても同じ (あるいは意味のある) 結果の得られる操作を Idempotent な操作という。Idempotent な操作では, 通信が失敗した場合などに, レスポンスを再送せず, クライアントが要求を再度行なえばよい。これに対応し VMTP はサーバ側で再送のためのレスポンスの保存を省略する。

2.5 通信のためのカーネルの最適化

V の IPC の基本はブロック型のメッセージ送信であるため, プロセスは通常一度に一つのメッセージ送信しか行なわない。カーネルではプロセス管理構造体中に VMTP ヘッダを含むことで, バッファ

管理を簡単化している。また、ヘッダの一部のフィールドはプロセス初期化時に値を設定してしまう。

2.6 マルチキャスト

V システムの IPC の特徴の一つにマルチキャストがある。V の通常の IPC は、通信のエンド・ポイントとしてプロセスを指定する。V ではプロセスの他にプロセス・グループを通信のエンド・ポイントとして指定できる。プロセスはプロセス・グループを生成、削除したり、参加、離脱することができる。プロセス・グループに対する Send はそのときプロセス・グループに参加しているプロセスにマルチキャストされ、最初の Reply メッセージを受け取ると待ちを解除される。2 個目以降の Reply は GetReply で受け取る。

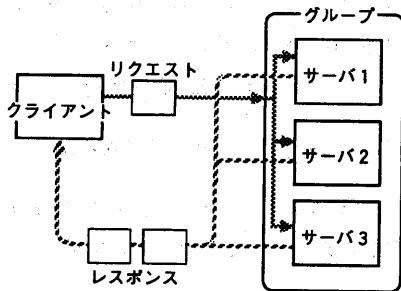


図 3: マルチキャスト

プロセス・グループに対するマルチキャストの典型的な使用例は名前管理である。V では専用のネームサーバを持たず、資源の各サーバが分散して名前を管理する。文字列により名前を管理する全サーバは、CSNH(Character String Name Handling) サーバグループと言う静的に定義されたプロセス・グループに属している。クライアントが文字列名を CSNH サーバグループにマルチキャストすると、その名前に対応する資源を管理しているサーバが Reply することで、サーバが特定される。指定した文字列名が存在しないことは、マルチキャストが(何回かのリトライの後で)タイムアウトすることで知ることができる。

2.7 カーネルサーバ

分散カーネルは IPC 以外にプロセス管理、メモリ管理、デバイスドライバなどの機能を提供する。これらは Send で要求を出すサーバの形で実現されている。プロセスのスケジューリングはプライオリティ・ベースの単純なもので、タイムスライスなどはカーネル外のスケジューラ・プロセスが実現する。マルチプロセッサに対してはプロセッサ毎にレディキューをもつスタティックバインドを採用している。

V のプロセスはライト・ウェイト・プロセスであり、複数のプロセスが 1 仮想空間を共有できる。空間を共有するプロセスをチームと呼ぶ。通常チームはプログラムに対応している。

3 サーバ

ファイルや認証管理などの機能はプロセスレベルのサーバとして実現されている。サーバ実現上共通性のある I/O とネーミングは各サーバに共通のプロトコルが規定されており、サーバの扱いを統一化すると共に、プロトコルを扱うサーバ側ライブラリを提供することでサーバの実現を容易にしている。ここでは I/O とネーミングのプロトコルの特徴を簡単に紹介し、サーバの例としてチームサーバ、Internet サーバ、V サーバについて述べる。

3.1 I/O プロトコル

V システムでは UIO と呼ばれる I/O プロトコルが規定されている。詳細は [2] を参照していただくとして、ここでは 2~3 の特徴を述べる。

アプリケーションに対しては getc, putc, fseek などの UNIX ライクなインタフェースが提供されている。このアプリケーション・インタフェースはクライアント側ライブラリによってクライアント-サーバ間プロトコルに対応付けられる。ライブラリ内ではキャッシュによってサーバへの通信回数を減らし、効率をあげている。また、ファイル中のリード/ライト・ポインタをクライアント側ライブラリで管理している。これによりリードはサー

バの状態を変えない(リードポインタが進まない) Idempotent なものになる。Idempotent な通信に対しては VMTP がリブライを再送用に保存しないという最適化を行なうため、効率のよい I/O が実現される。このようにできる限り多くの処理をクライアント側ライブラリで行なうことにより、アプリケーション・プログラマには効率化の詳細を意識させずに、共有資源であるサーバへの負荷の集中を避けている。

3.2 ネーミングプロトコル

V システムのネーミングは資源を管理する各サーバが資源の名前も管理する分散管理方式である [7]。これはネームサーバによる集中管理よりも効率的であり耐故障性も高い。2.6で述べたように、資源を文字列名で管理する各サーバは CSNH サーバグループに属している。クライアントは文字列名を CSNH サーバグループにマルチキャストすることでサーバを特定するが、通信の回数を減らすために名前前のプレフィクスとサーバの対応のキャッシュを持っている。サーバがクラッシュして再起動された場合にはサーバの pid が変わりキャッシュの一貫性が崩れる。これに対しては、とりあえずキャッシュの値を使用して通信してみて、値が古くなっていたらマルチキャストして正しいサーバを探すと同時にキャッシュを更新する。このような考え方を問題指向共有メモリ [1] と呼んでいる。

ネームスペースはシステムグローバルな木構造になっている。

3.3 チームサーバ

チームサーバはチーム(プログラム)のロード、実行、モニタなどの制御を行なう。チームサーバには次のような機能がある。

- チームのロード、プライオリティの設定
- チームの終了
- チームのラウンド・ロビン・スケジューリング

- デフォルト例外ハンドラの提供。(実際には例外を起こしたプロセスに対しデバッガを起動する。)
- プロセスのリモート実行の際のホスト選択のための負荷情報の管理

チームのロード要求はリモートなチームサーバに対してもローカルな場合と同様に行なえる。これによりネットワーク透過なプログラムのリモート起動が実現される。チームロードのクライアント側ライブラリでは、その時点での各ホストの負荷をチームサーバに問い合わせ、もっとも負荷の低いサーバにチームをロードするという機能が実現されている。

3.4 Internet サーバ

V システムはホスト間プロトコルとして VMTP を採用している。Internet サーバは VMTP が実装されていない(V以外のOSが動いている)ホストとの通信を行なうサーバである。このサーバは

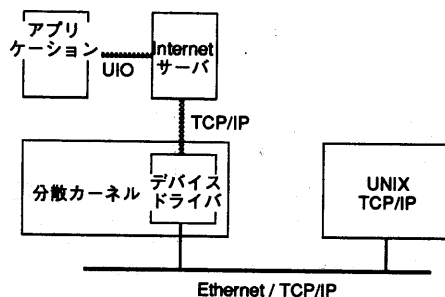


図 4: Internet サーバ

基本的にはプロトコル変換機であり、VのアプリケーションからはVのI/Oプロトコルでアクセスできる。Internetサーバはこれを指定されたプロトコルに変換し、ネットワーク・デバイスと直接入出力することで、VMTPの実装されていないホストと通信する。現在プロトコルとしてはTCP/IPがサポートされているが、Internetサーバの構造はプロトコル独立なフレームワークとプロトコル

依存部分とが分離されており、新たなプロトコルの追加が容易になるように考慮されている。

3.5 V サーバ

V サーバは他の OS(現在は UNIX のみサポートされている)のファイルシステムやプロセスを、V のファイルやプロセスのように扱うためのサーバである。V システムには、V の I/O プロトコルでアクセスする独自のファイルサーバがある。システム中には複数のファイルサーバが存在できるようになっている。V サーバは他の OS 上で動作するサーバで、V のネーミング・プロトコルと I/O プロトコルを認識し、それを他の OS のファイル I/O に変換する。これにより、V のアプリケーションからは V のファイルサーバをアクセスしているのと同じプロトコルで実際には他の OS のファイルがアクセスできる。V のプロセッサ間プロトコルは VMTP なので、V サーバが動作するマシンは VMTP が実装されている必要がある。

また、V のプロセスとの間にパイプを作り、他の OS のプログラムを他の OS 上のプロセスとして実行することもできる。

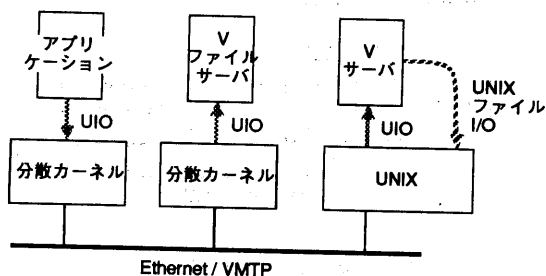


図 5: V サーバ

3.6 多重化ファイルサーバ

著者の V プロジェクトでの研究テーマであった多重化ファイルサーバについて簡単に述べる。

ネットワーク上でファイルサーバを共有し、ファイルが数ヶ所のホストに集中すると、ファイルサーバ

の障害がシステム全体に及ぼす被害は大きくなる。そこで、ファイルサーバの耐故障性が重要な課題になる。

分散システムでは、同じファイルを持つファイルサーバを異なるホストで複数動作させ、任意のサーバからリードすることで、耐故障性が得られる。分散システムでは一般的に、故障に対するホストごとの独立性が高いので、ディスクの二重化などに比べより高い耐故障性が実現できる可能性がある。

多重化サーバのデータを同時に更新するために、UIO のリード/ライト・プロトコルを“多重系全サーバヘライト/任意サーバからリード”に拡張した。多重系を構成するサーバは同じプロセスグループに属するようにして、クライアントがそのグループにライト要求をマルチキャストすることで、全サーバライトは効率よく実現できる。リードはオープン時にグループに名前参照要求をマルチキャストし、最初にリプライを返したサーバのファイルをオープンすることで任意サーバを選択した。これはオープン時に最も負荷の低いサーバが選ばれることになり、リード処理の負荷を動的にバランスさせるという副作用も持っている。リードエラーの場合、再度別のサーバをオープンし、同じブロックを再リードすることで、処理を続行する。この時、リード操作の Idempotency がファイルサーバにまたがって維持されるため、サーバの切り替えは意識されない。以上のような I/O プロトコルの拡張は全てライブラリ内に閉じており、アプリケーションプログラムには意識させずにファイルの多重化が実現できた。

拡張したプロトコルのオーバーヘッドをプロトタイプを作成して評価したところ、多重度 1~3 の範囲ではリード/ライト/リードオープン処理にはほとんどオーバーヘッドがないことがわかった。ライトオープンでは多重度 2 で 2.6 倍、3 で 2.8 倍程度の時間がかかった。イーサネットでのマルチキャストが物理的にも 1 パケットで実現されていることが、低いオーバーヘッドに寄与していると思われた。

耐故障サーバの実現のためには、復旧したサーバのデータの同期、ネットワークのパーティションへの対応などの問題が残されている。

4 ユーザ環境

V システムは全体としてはワークステーション群を一つの TSS マシンに見せる分散 OS と見ることができ、ブート直後やユーザが使用していない状態のワークステーションでは、分散カーネル上にチームサーバ、例外サーバ、exec サーバ、ウィンドウ管理サーバが動いている [6]。

4.1 Exec サーバ

Exec サーバは UNIX の shell 風のコマンド・インタプリタで、ヒストリ、alias、入出力のリダイレクト、パイプなどの機能やコマンドの 1 ラインエディットの機能がある。

コマンドのサーチ順序は環境変数 PATH で指定できる。コマンドがこのサーチパスにない場合には、カレント・ディレクトリを提供するサーバでリモート実行される。カレント・ディレクトリが V サーバが提供する UNIX ファイルの場合、まず V のコマンドが探され、次に UNIX コマンドが探されリモート実行される。つまりユーザから見ると V のコマンドも UNIX コマンドも同様に利用できることになる。最近では V 上の UNIX バイナリのシミュレータが開発され、UNIX コマンドは V 上で直接動作するようになったようである [8]。

通常、コマンドはログインしたホストで実行されるが、コマンドライン中にホストを指定することで、コマンドのリモート実行が行なえる。さらに、ホスト名として any を指定すると、その時点で誰もログインしていなくて、メモリと CPU の使用量の一番低いものがホストとして選択され、リモート実行される。

現在 V システムは Sun2, Sun3, μ VAX, DEC3100 などのアーキテクチャをサポートしている。V コマンドは CPU 毎に、Sun なら .m68k, VAX は .vax といった別のサフィックスが付けられている。Exec

サーバはコマンドを実行するホストによって自動的にサフィックスを選択するので、ユーザはホストの CPU タイプを意識する必要がない。もちろん any を指定して負荷の軽い任意のホストを選ぶ時にも自動的に適切なバイナリが選択される。

4.2 コマンド

現在 V システムでは 200 個程度のコマンドが提供されている。このうち約 100 個は UNIX から移植されたもので、cc, CC, emacs, mail, rcs, TeX, LaTeX などが含まれる。V サーバによって UNIX を V のファイルサーバとして使っている場合、UNIX コマンドもリモート実行して利用できるが、V コマンドに移植することで、コマンドが V のワークステーションでローカルに実行でき、UNIX マシンの負荷が分散できるため、使用頻度の高いものは積極的に V に移植されている。分散機能を生かしたコマンドとしてはパラレル make や実験的なものとして並列 smalltalk, 並列シミュレータなどが開発されていた。

4.3 OS の開発環境

V システムはファイル管理、プログラムロードなど通常の OS 機能がプロセスレベルのサーバとして実現されている。このため、このようなモジュールの開発も、対象のサーバをデバッガで対話的にデバッグしたり、システム全体を動作させたままで、再コンパイルして置き換えたりが可能であり、システムプログラムの開発効率が高い。また、分散カーネルにはデバッガが組み込まれており、カーネル内で例外が起こると自動的にこのデバッガが起動される。

V プロジェクトでは、開発環境として使用している安定したシステムと、実験中の不安定なシステムとを同じネットワーク中に接続しているが、バケット中の VMTP のバージョン番号を変えることで不安定なシステムが開発環境に干渉を与えないようにしている。

5 おわりに

V システムは 1986 年に V6.0 がリリースされた後も、積極的な拡張、書き直しがされている。現在では V6.0 とはかなり機能や内部構造が変化しているが、現在のところ次版のリリースの予定は定かではない。

また、ここに触れなかった話題として

- 共有メモリマルチプロセッサへの対応 [5].
- ファイル・キャッシュ・プロトコル [11].
- 分散アトミック・トランザクション・プロトコル.

などがある。いずれも研究開発中またはある程度成果がまとまった段階である。

参考文献

- [1] David R. Cheriton. Problem-oriented shared memory: A decentralized approach to distributed system design. In *the 6th International conference on Distributed Computer Systems*, pages 190–197. IEEE Computer Society, May 1986.
- [2] David R. Cheriton. UIO: A uniform I/O interface for distributed systems. *ACM Transactions on Computer Systems*, 5(1):12–46, February 1987.
- [3] David R. Cheriton. The V distributed system. *Communications of the ACM*, 31(3):314–333, March 1988.
- [4] David R. Cheriton. Versatile message transaction protocol (vmtp). Technical Report RFC 1045, SRI Network Information Center, February 1988.
- [5] David R. Cheriton, Hendrik A. Goosen, and Patrick D. Boyle. Multi-level shared caching techniques for scalability in vmpmc. In *the 16th Annual International Symposium on Computer Architecture*, May 1989.
- [6] David R. Cheriton and Keith A. Lantz. *V-System 6.0 Reference Manual*. Stanford University, July 1986.
- [7] David R. Cheriton and Timothy P. Mann. Decentralizing a global naming service for improved performance and fault tolerance. *ACM Transactions on Computer Systems*, 7(2):147–183, May 1989.
- [8] David R. Cheriton, Gregory R. Whitehead, and Edward W. Szynter. Binary emulation of Unix using the V kernel. In *USENIX Summer Conference*, pages 73–85. USENIX, June 1990.
- [9] David R. Cheriton and Carey L. Williamson. Network measurement of the VMTP request-response protocol in the V distributed system. In *SIGMETRICS 87*. ACM, 1987.
- [10] David R. Cheriton and Carey L. Williamson. VMTP as the transport layer for high-performance distributed systems. *IEEE Communications*, pages 37–44, June 1989.
- [11] Cary G. Gray and David R. Cheriton. Leases: An efficient fault-tolerant mechanism for distributed file cache consistency. In *the 12th ACM Symposium on Operating System Principles*, pages 202–210. ACM SIGOPS, December 1989.