

第13回 仮想計算機・プロセス間通信

前回のスライド

重要: モノリシックカーネルとマイクロカーネル

マイクロカーネル

モノリシックカーネル

OSサーバ

OS(カーネル)

マイクロカーネル

OS(カーネル)

OSを複数のモジュールに分割

構成: OS全体を1モジュールで構成

・特権モードで動作するマイクロカーネルにはプロセス管理, メモリ管理, プロセス間通信等OSの最低限の機能のみを残す

利点: 効率が良い(オーバーヘッドが少ない)

・他の付加的機能は, 非特権モードで動作するように階層化する

欠点: OSの肥大化により, 管理が困難(機能の追加・変更が難しい)

・特権モード(注)で動作するプログラム量が多くデバッグが難しい(特権命令を使用)

利点: 管理が簡単, デバッグしやすい

利点: 管理が簡単, デバッグしやすい

欠点: オーバーヘッドが大きい

注: 特権モード(カーネルモード)

注: 非特権モード(ユーザーモード)

前回のスライド

Windows NT/2000/XP

WIN32 アプリケーション

DOS アプリケーション

WIN16 アプリケーション

POSIX(注) アプリケーション

OS/2 アプリケーション

ログオンプロセス

Windowsサブシステム(非特権モードで動作)

WOW

POSIXサブシステム

OS/2サブシステム

セキュリティサブシステム

システムサービスAPI

NT Executive(特権モードで動作)マイクロカーネル

NTカーネル

Hardware Abstraction Layer

ハードウェア

注: POSIX (Portable Operating System Interface: Unixの標準の一つ)

前回のスライド

(重要) Windows NTクライアント-サーバ構造

アプリケーション

サブシステム

NT Executive

kernel

Win32クライアント

Win32サーバ

OS/2クライアント

OS/2サーバ

POSIXクライアント

POSIXサーバ

マイクロカーネルでは, アプリケーション(利用者プログラム, クライアント)は, カーネルを経由したメッセージ通信を行うことによりOSのサブシステム(サーバ)の機能を利用する(クライアント, サーバともに利用者モードで動作する)

クライアント, サーバという用語はプロセスの機能分担を指す。コンピュータ間の関係ではない。(クライアント≠PC等の装置, サーバ≠ワークステーション等の装置)

前回のスライド

重要: マイクロカーネルのシステムコール

Windows NTの場合

クライアントプロセス Win32 APL

サーバプロセス Windowsサブシステム

①処理を依頼(システムコール)

②サービス依頼(システムコール)

③結果通知(システムコール)

非特権モード

特権モード

LPC Local Procedure Call

LPC機構

NT Executive

図の例では, クライアントプロセスは, Windowsサブシステム(サーバプロセス)に対して, 1回のシステムコールを実行。しかし, 実際には, サーバプロセス(非特権モードで動作)もカーネルであるNT Executiveに対し, 2回のシステムコールを実行している(計3回(注))。

注: モノリシックカーネルでは1回のみ。

マイクロカーネルの欠点: オーバーヘッドが大きい(割り込み分析, モード切替, レジスタの退避など種々の処理が必要)。

重要: 仮想計算機 (VM: Virtual Machine)

コンピュータのハードウェア資源を仮想的に分割し, 複数の独立した実行環境を提供する機構(注)。仮想計算機(VM)毎にOSをローディングして実行できる。OSは同一でも異なっても良い。

・仮想計算機モニタがVM毎のハードウェアをエミュレートし, 実ハードウェアと同一のインタフェースを提供。

・CPUスケジューリングと仮想記憶により, 各VM専用の仮想CPUと仮想メモリを提供。

・ディスクを領域分割し, VM毎に独立したファイル管理などができるようにする。

利点: 強力な保護機能を持つ(他の仮想計算機とは, 完全に独立)

欠点: 命令の解釈実行などのオーバーヘッドがある(VMあたりのCPU性能<1/n)

仮想計算機間で, 直接の資源共用ができない(ネットワーク経由の通信などで共用)

注: 当初は, OSのデバッグや乗換えを効率化する目的で開発された。

1台のハードウェア上にn台のコンピュータ(VM)

VM1 VM2 VM3

AP OS OS OS

仮想計算機モニタ(ハイパーバイザ)

ハードウェア

1台のハードウェア上に1台のコンピュータ

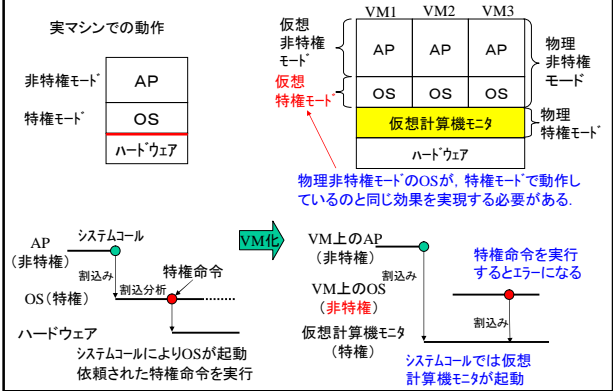
AP OS ハードウェア

同一仮想ハードウェアインタフェース

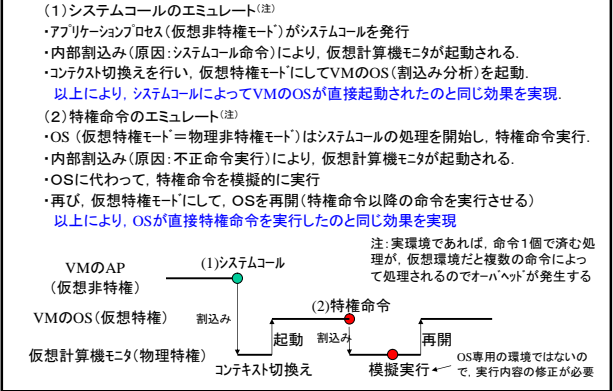
実ハードウェアインタフェース

エミュレート(Emulate): まねる, 模倣する

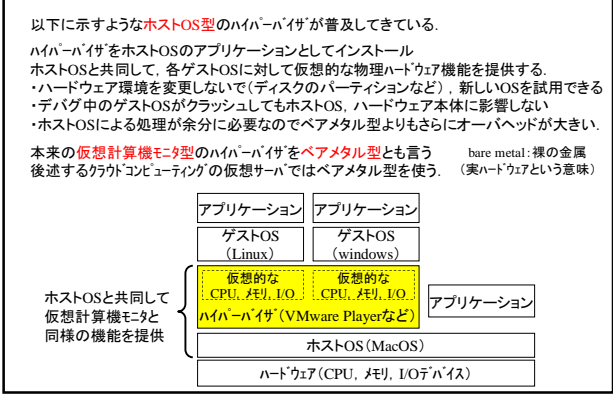
重要: 仮想計算機の実現



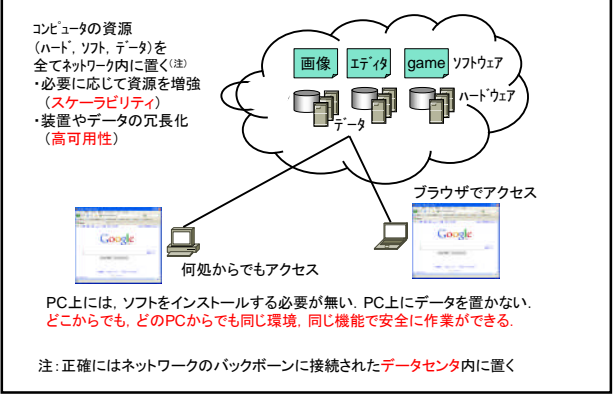
重要: 仮想計算機の動作例



ホストOS型ハイパーバイザによる仮想化



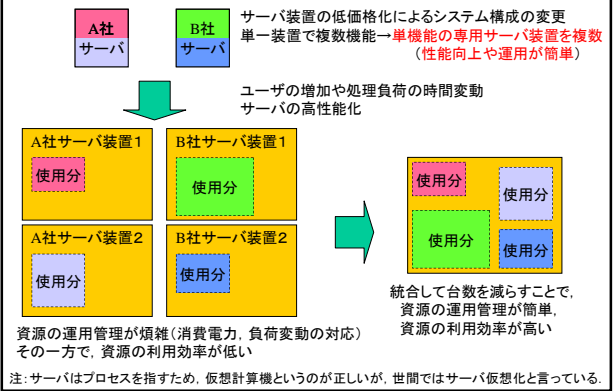
クラウドコンピューティングの概念



クラウドコンピューティングの要素

- ハウジング
 - 高速なネットワークに接続されたデータセンタに自分のサーバを預かってもらい、保守や管理をアウトソースする。
- ホスティング
 - データセンタ内のサーバおよびソフトウェアを借りて使う。
 - 特に、複数のユーザでwebサーバの機能をシェアして使う。
- SaaS (software as a service)
- ASP (アプリケーション・サービス・プロバイダ)
- SOA (サービス指向アーキテクチャ)
 - アプリケーションプログラムを借りて使う
- サーバの仮想計算機化 (仮想サーバ、仮想マシン)
 - 1台のサーバを複数のサーバとして各ユーザに提供
- グリッドコンピューティング
 - 複数のコンピュータを1台のコンピュータのように見せる

サーバ仮想化^(注)の必要性



プライベートクラウドの提供とその運用

プライベートクラウド: 自社専用のクラウド
最近ではデータセンタの設備を共有してプライベートクラウドを提供するソリューションがトレンドとなっている。

必要な性能を維持するため、物理サーバ装置を移動

物理サーバ装置1 物理サーバ装置2

仮想サーバ 仮想サーバ 仮想サーバ 仮想サーバ 仮想サーバ 仮想サーバ

仮想スイッチ 仮想スイッチ 仮想スイッチ 仮想スイッチ 仮想スイッチ 仮想スイッチ

レイヤ2スイッチ (VLAN switch)

仮想サーバの移動に合わせて VLAN の設定を変更

性能向上の手法

スケールアップ型
×故障発生時の影響範囲大
×高価なハードが必要
×ユーザ増加への対応が困難

2000万ユーザ収容

スケールアウト型(クラウド化)
○故障発生時の影響範囲小
○比較的安価なハードで構成可能
○ユーザ増加への対応が容易
但し、ソフトウェア技術が必要

200万ユーザ収容

20万ユーザ収容

装置能力を拡大

装置台数を増やす

クラウドコンピューティングプラットフォーム
Google App Engine, Amazon EC2, Windows Azure...

プロセッサ性能と価格

コスト比 1:5

性能比 1:2

プロセッサの性能向上も マルチコア化が潮流

排他制御の実現方法とプロセス間通信

第3回のスライド

P1のプログラム P2のプログラム

共有変数

①システムコール ②チェックオン ③割り込み ④システムコール ⑤待機状態に ⑥再開

⑦システムコール ⑧オフ ⑨レディ状態に ⑩終了

問題点1の解決法: システムコールにより、OSがフラグを操作する②④⑧
なお、OSがフラグを操作中は、割り込み禁止(またはフラグのテストとセットを一度に行う命令を使う)

問題点2の解決法⑤⑨⑩
フラグオンの場合、プロセス(P2)を待機状態にし、P1のクリティカルセクション終了後にレディ状態にする(左図のようなプロセス間通信とみなせる)

プロセス間通信

- 協調が必要なプロセス間でデータの授受(通信)を行う
 - 同期問題はプロセス間通信の一部(危険区域を出た通知)
- 共有メモリ方式(APがプログラムの中で通信機能を実現)
 - 複数のプロセスがアクセス可能な共有メモリを用意し、この領域を利用してデータの受け渡しを実現する
- メッセージシステム方式(OSが提供する通信機能を利用)
 - システムコールを用いて、メッセージの形で通信する
 - send(destination, message)/receive(source, message)
 - 直接通信と間接通信

プロセス間通信(共有メモリ方式)

共有メモリ方式(APがプログラムの中で通信機能を実現)
複数のプロセスがアクセス可能な共有メモリを用意し、この領域を利用してデータの受け渡しを実現する

制約バッファ問題の例

生産者 共有メモリ 消費者

buffer[i]=nextp
i=i+1%n

buffer[0]
buffer[1]
buffer[2]

nextc=buffer[j]
j=j+1%n

buffer[0]~buffer[n]は、生産者、消費者が共有するメモリ

lock/unlock, セマフォなどを用いた排他制御も必要であり、APが複雑となるメモリへのアクセスのみで通信が可能のため、効率が良い

プロセス間通信(メッセージシステム方式:直接通信)

OSが提供する通信機能を利用: システムコールを用いて、メッセージの形で通信する

直接通信では、通信相手のプロセス名を指定して、メッセージを送受する
通信相手が固定的に決まっている1:1通信を簡単に実現
モジュール性が悪い: プロセス名を変更すると相手プロセスに影響する

send(P2, message) P1 P2 receive(P1, message)

OS

```
producer: /*生産者*/
do{
...
データを読んでnextpへ
send(consumer, nextp)
} while (1)

consumer: /*消費者*/
do{
receive(producer, nextc)
...
nextcのデータを出力
...
} while (1)
```

プロセス間通信(メッセージシステム方式:間接通信)

send(A, message)send(A, message)receive(A, message)

P1P2P3

A

OS

間接通信では、チャネル(channel)^(注)を介してメッセージを送受する
チャネルを共有するプロセス間にリンクが作られ、通信ができる
1:多通信(受け手が複数)、多:1通信(送り手が複数)、多対多通信(左記の組み合わせ)が可能

クライアント〜サーバ間の通信(多対1通信)を実現
メッセージの受信(サーバ)は、通信相手(送信プロセス名)を指定しない方式
(チャネルの変数名を指定し、受信時に送信プロセスのidを得る)
直接通信方式よりもチャネルを介した間接通信方式の方が柔軟である。

注: ポート(port)、メールボックス(mailbox)とも言う
TCP、UDPで用いるポート番号は、チャネル識別子である

参考:同期式通信と非同期式通信

同期式通信

非同期式通信

受信プロセスがメッセージの受信処理を終えるまで待たされる(待機状態となる)
受信処理完了により、send命令が終了。
(メッセージが届いたことが分かる)

受信プロセスが受信しなくてもメッセージがバッファに実行を続けることができる。
但し、send命令で、メッセージが届いたかどうかは分からない

通常のクライアント・サーバの通信はこの形式

短時間に処理を終わらせる必要がある場合は、この形式が採られる。

プロセス間通信の構造化

5プロセス 8メッセージ

P1P2P3P4P5

構造化されていない通信

P1P2P3P4P5

構造化された通信

各プロセスの機能分担を見直し、見通しを良くする

重要:クライアントサーバモデル

プロセスを2種類に分ける

クライアント

サーバ

サーバ

クライアント1クライアント2クライアント3

クライアントサーバの通信パターン

クライアント

サーバ

while(1){ send(); receive(); }

while(1){ receive(); send(); }

クライアントは、要求メッセージを1回送信(send)し、応答メッセージを1回受信(receive)
サーバは、要求メッセージを1回受信(receive)し、応答メッセージを1回送信(send)
双方、このパターンを守り、必要な回数繰り返す。

① ②

第9回のスライド

Windows NT/2000/XP

WIN32 アプリケーションDOS アプリケーションWIN16 アプリケーションPOSIX^(注) アプリケーションOS/2 アプリケーションログオンプロセス

サブシステム(非特権モードで動作)Windows サブシステム(Kernel)MS-DOSサブシステムPOSIX サブシステムOS/2 サブシステムセキュリティサブシステム

システムサービスAPINT Executive(特権モードで動作)マイクローネルNTカーネルハードウェア

注: POSIX (Portable Operating System Interface: Unixの標準の一つ)

