

Linux カーネル入門

VA Linux Systems ジャパン(株) OSDN 事業部
安井 卓 <tach@valinux.co.jp>

2002.10.9

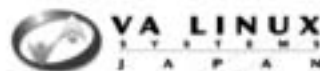
Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



もくじ

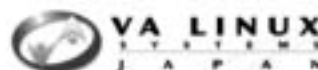
- Linux カーネルとは
- カーネル主要機能概要
- そのほかのトピック
- 参考文献・ポインタ

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



Linux カーネルとは

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



Linux システム



様々なツールやライブラリと一緒に「Linux」として配布しているが、本当に Linux といえるのはカーネルのみ

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



Linux システムでのカーネルの役割



ハードウェアデバイス



VA LINUX
SYSTEMS
J A P A N

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.

カーネル(というソフト)の特徴

- ・ハードウェアの制御を行う
 - ・デバイスドライバ
 - ・割り込み処理
- ・アプリケーション(ユーザプログラム)に実行環境を提供する
 - ・システムコール
 - ・メモリの管理
- ・カーネルは自発的になにもしない
 - ・システムコールや割り込みなど、イベント駆動型のソフトウェア

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.

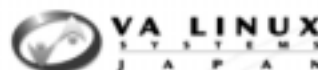


VA LINUX
SYSTEMS
J A P A N

Linux カーネルの特徴

- POSIX 準拠の OS である (UNIX 互換)
- モノリシックカーネルである
- カーネルモジュールの動的なロードができる
- オープンソース (GPL) である
- 多くのアーキテクチャに対応している
- 開発速度が速い

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



Linux カーネルの主要機能

システムコール インターフェイス

時間管理

プロセス管理
(スケジューラ)

空間管理

割り込み管理

ネットワーク

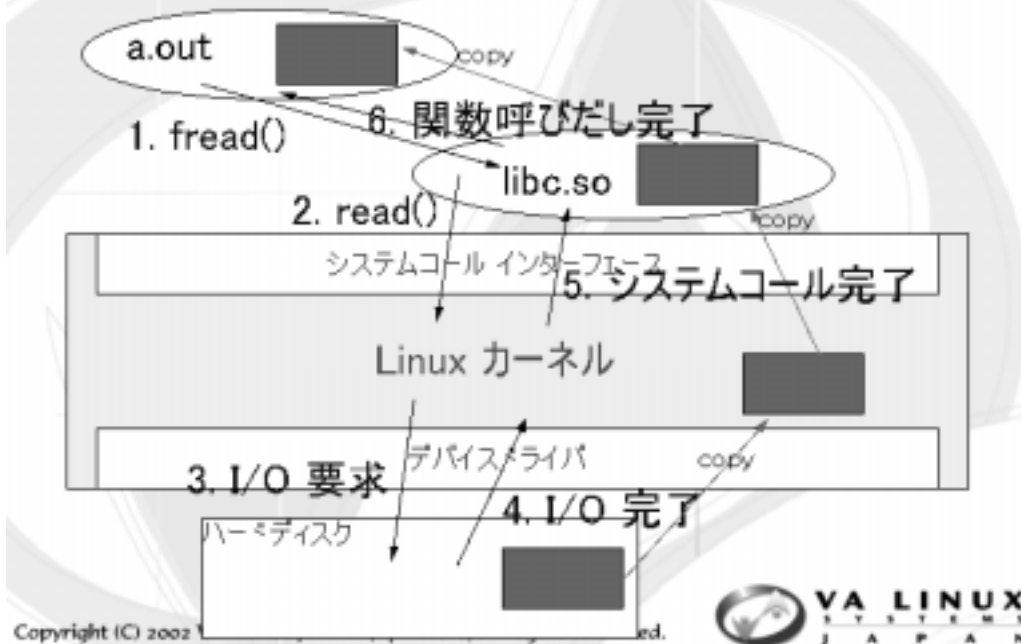
ファイルシステム

ブロックデバイスドライバ | キャラクタデバイスドライバ | ネットワークデバイスドライバ

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



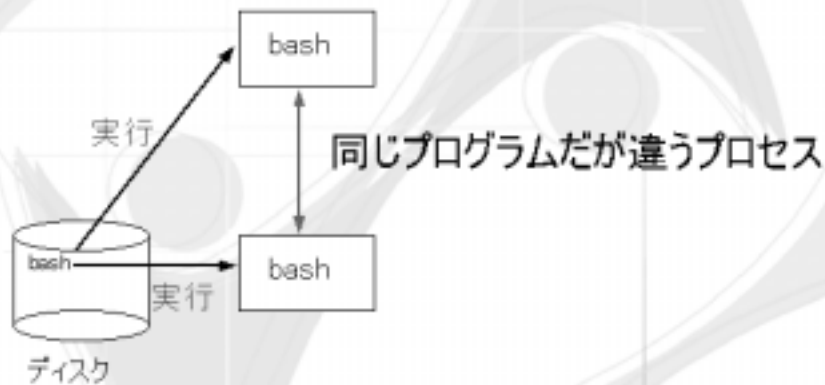
Linux カーネル処理例



Linux カーネル 機能概要

プログラムとプロセス

プロセス … プログラムが動いている状態



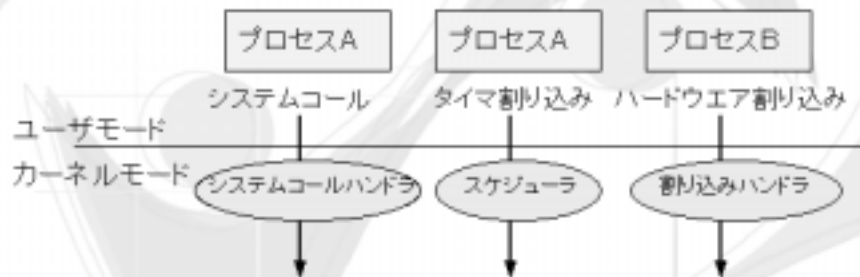
Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



カーネルモード・ユーザモード

カーネルモード = 特権モード

- ・すべてのカーネルデータ・コードにアクセスできる
- ・プロセスはふだんはユーザモードで動く

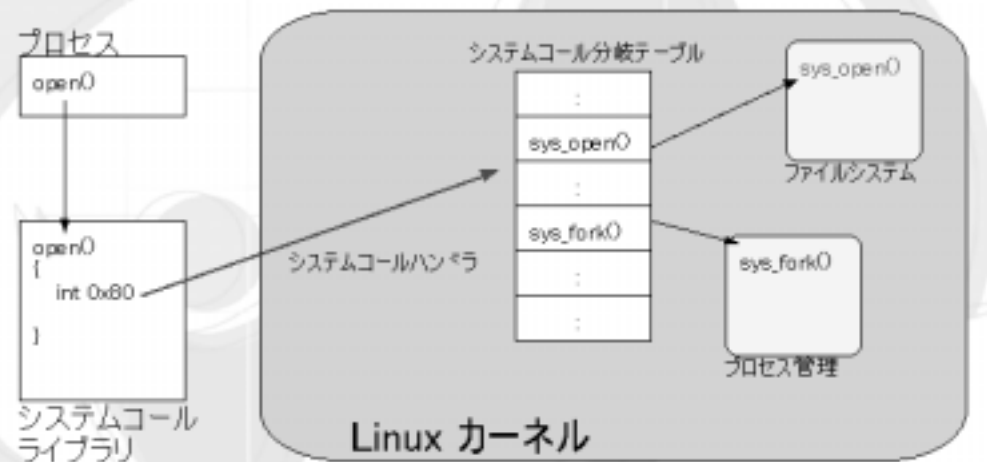


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



システムコール

- Linux カーネルに対するソフトウェアインターフェース
- カーネルに対するソフトウェア割り込み



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



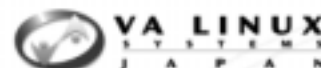
プロセスの情報

各プロセス毎にプロセス情報構造体におさめられている

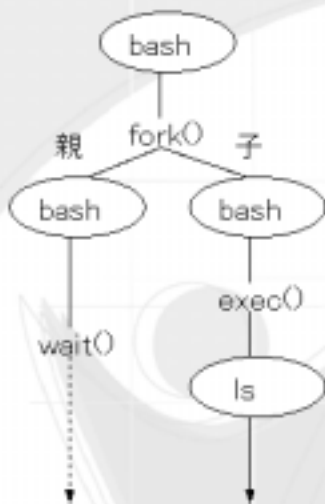
- プロセスの識別子
- プロセス名
- プロセスの状態
- オープンしているファイルの情報(ファイルディスクリプタ)
- メモリ情報
- 優先度(プライオリティ)
- コンテキスト保存領域
- シグナルハンドラ
- …など



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセスの生成



1.fork()システムコールで自分自身(親プロセス)の複製を作成

2.(子)exec()システムコールで、自分自身をlsコマンドで置き換え

3.(親)wait()システムコールで、子の終了を待つ

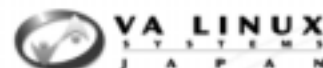
Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



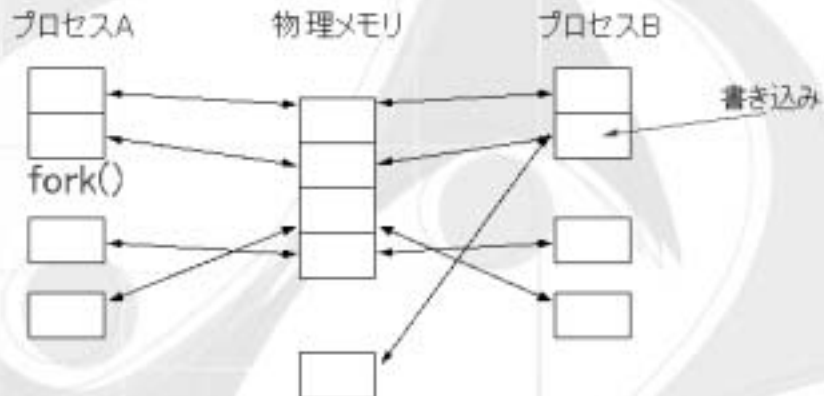
fork システムコール

- ・自分自身の複製プロセスを生成するシステムコール
- ・プロセスを生成する唯一の手段
- ・親プロセスの持つすべての資源も複製する
 - ・ファイルディスクリプタ
 - ・シグナルハンドラ
 - ・プロセス空間(コピーオンライト属性)
- ・プロセスIDはカーネルが適当に割り当てる

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



コピーオンライト



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



exec システムコール

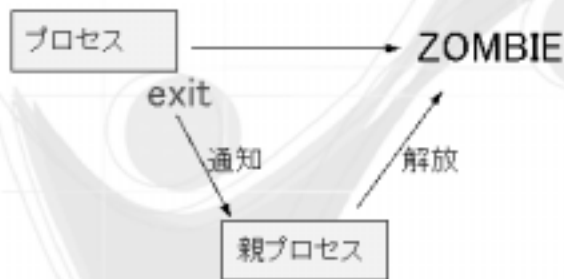
- ・自分自身を新しいプロセスに置き換える
- ・すべてのプロセス空間を解放し、新しいプロセス空間を設定
- ・ファイルディスクリプタなどの管理情報は引き継ぐ

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



exit システムコール

- ・プロセスの終了時に呼び出す
- ・すべてのプロセス空間を解放
- ・ファイルディスクリプタなどの管理情報は残す



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセス管理

LinuxはマルチタスクOSなので、(CPUがひとつでも)複数のプロセスが同時に動く。

…というより、動いているように見える。
実際は時分割で複数のプロセスを切り替えながら動く。
(プロセスディスパッチ)

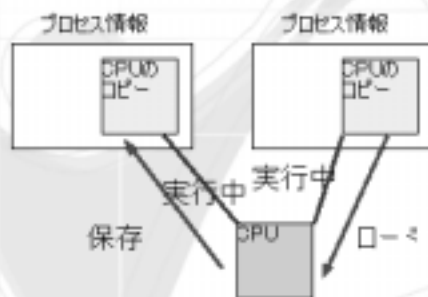


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセスディスパッチ(切り替え)

- ・現在動作しているプロセスから、次に動作すべきプロセスに切り替える動作

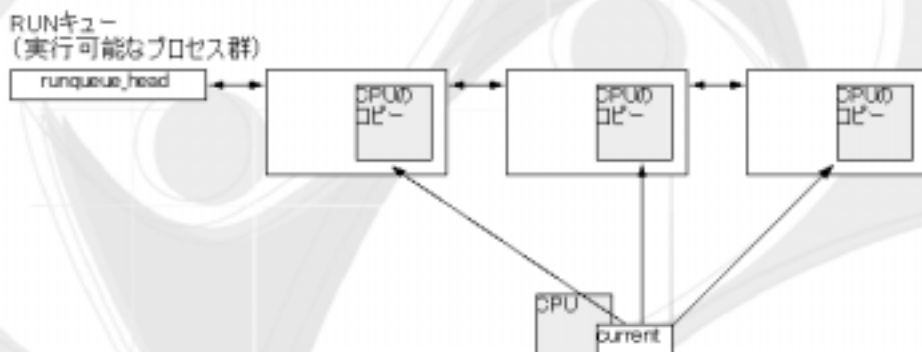


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセススケジューリング

- ・実際にどのプロセスを動かすかの選択
- ・プロセススケジューラ



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセススケジューラ



プロセスの優先度を計算

各プロセスに対して、CPUの1回の持ち時間を与える
(優先度に応じて、というか優先度そのもの)

優先度が高いものから順番に実行

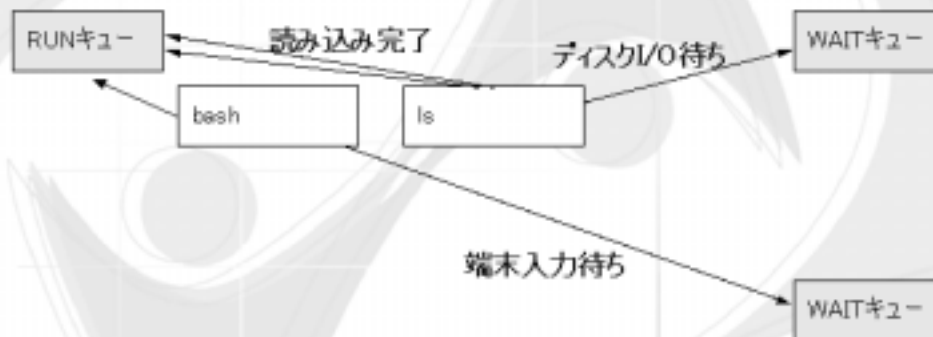
全プロセスの持ち時間がなくなったら優先度再計算
以下繰り返し...

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセスの待ち合わせ

ほとんどのプロセスは待ち状態にある

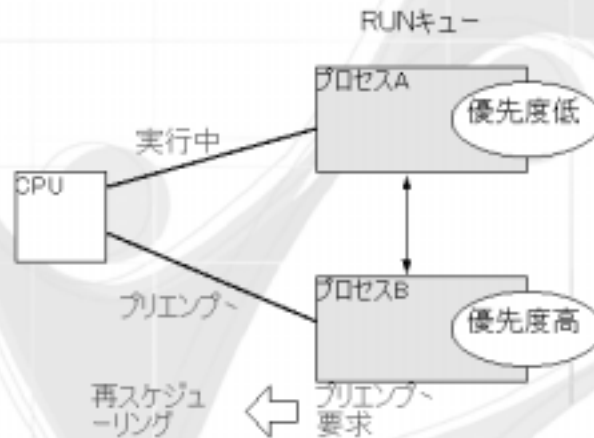


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プリエンプション

- あるプロセスが、カレントプロセスからCPU実行権を奪い取る処理



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.

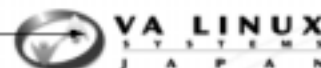


プライオリティ(優先度)

- 固定プライオリティ
 - niceコマンドなどで与えられる優先度
- 変動プライオリティ
 - 実行時間により動的に変更される
 - 実行時間が長いものほど優先度が下がる
 - I/O処理などが多いプログラムは優先度高 → 応答性がよくなる
 - 優先度がそのままCPU時間に



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセスとスレッド

- clone システムコールで生成
- ・実行単位はプロセスと同じ
 - ・資源のコピーを行わずに共有する

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



プロセスの状態

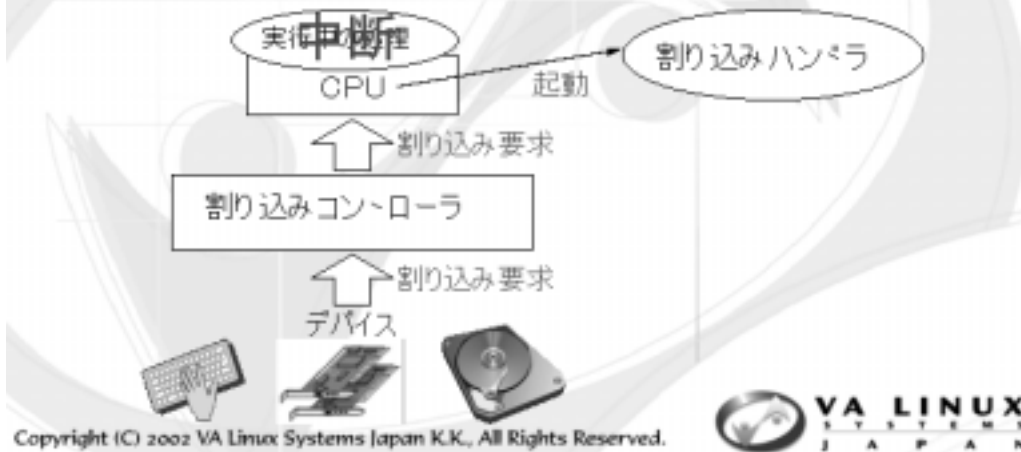


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



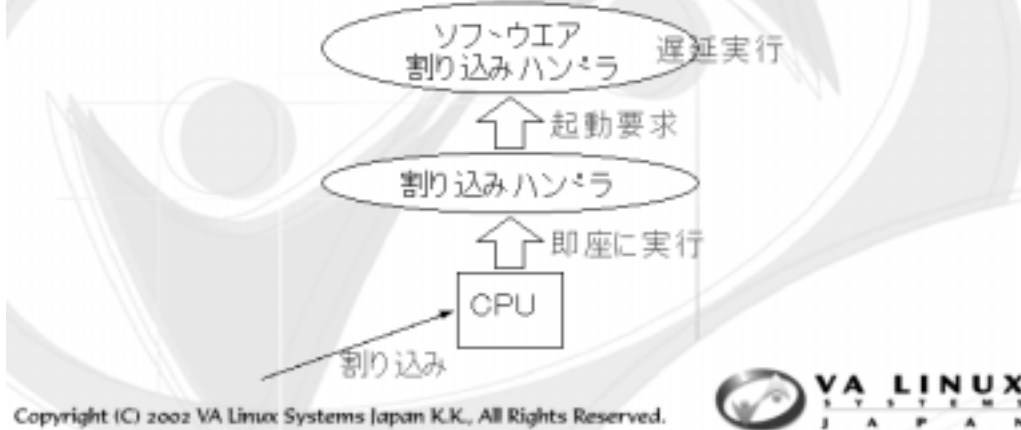
割り込み

- ハードウェアで発生したイベントを拾ってデバイスドライバを起動するためのしくみ

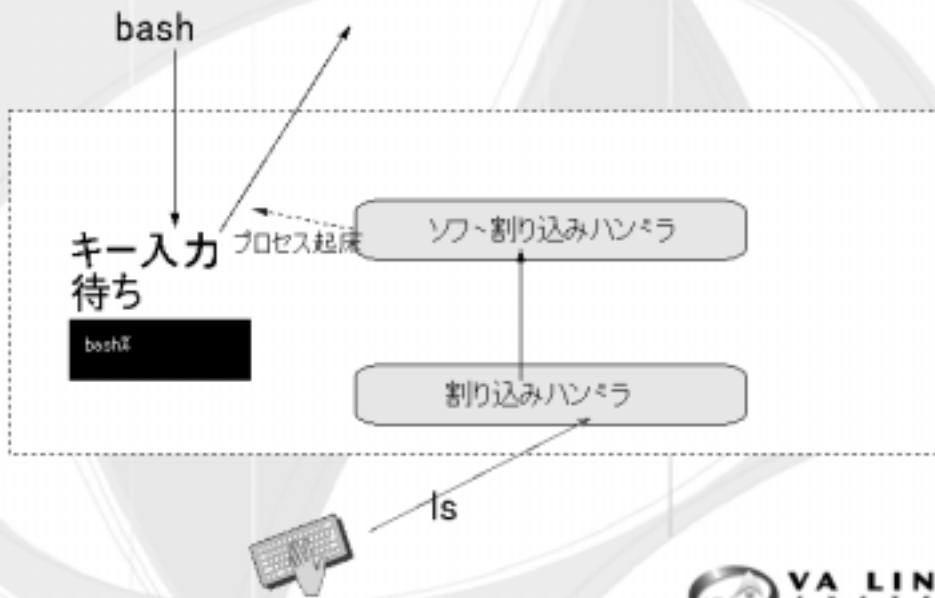


割り込みハンドラ

- 2つのハンドラ
 - 割り込みハンドラ
 - ソフトウェア割り込みハンドラ(BHハンドラ)



典型的なI/O割り込み処理例

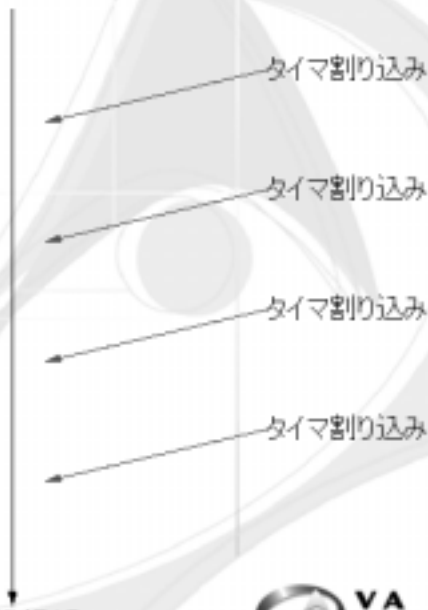


Copyright (C) 2002 VA Linux Systems Japan K.K. All Rights Reserved.



時間管理

- ・時刻を正確に刻む
- ・次元処理の実行
- ・定期的なタイマ割り込み



Copyright (C) 2002 VA Linux Systems Japan K.K. All Rights Reserved.



タイマ割り込みハンドラ

・タイマ割り込みハンドラ

- ・jiffies++
- ・2段目の時計 (BHハンドラ) の起動要求

BHハンドラ

- ・カレンダー (GMT時刻) の更新
- ・ロードアベレージの計算
- ・タイマリスト (時限処理メカニズム) の実行



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



ファイルシステム

- ・論理的なファイル構造と物理的なディスクブロックの対応を管理
- ・内部で階層化されている



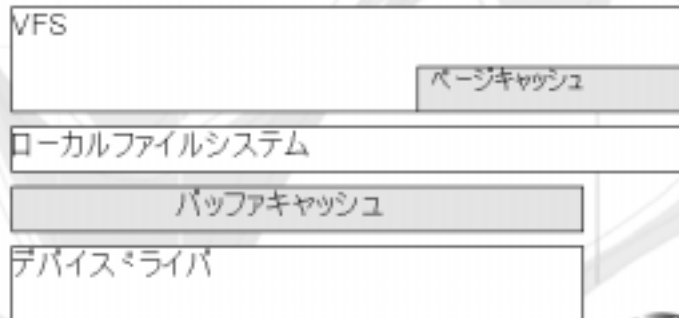
Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



ページキャッシュとバッファキャッシュ

ファイルアクセスを高速化するためのキャッシュ

- ・ページキャッシュ
 - ・ファイルそのもののキャッシュ
- ・バッファキャッシュ
 - ・ディスクブロックのキャッシュ

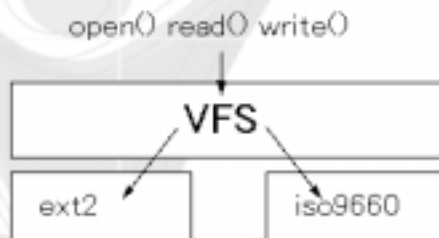


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



仮想ファイルシステム (VFS)

- ・共通ファイルシステムレイヤ
- ・ファイルシステム間の違いを意識せずにファイルにアクセス可能
- ・ファイルが物理的にどう配置されているかは関知せず、論理構造とインターフェースを定義



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



共通ファイルモデル

- VFS がサポートするファイルシステムモデル
伝統的UNIXファイルシステムを反映している
- スーパーブロックオブジェクト
- iノードオブジェクト
- ファイルオブジェクト
- dエントリオブジェクト

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



ローカルファイルシステム

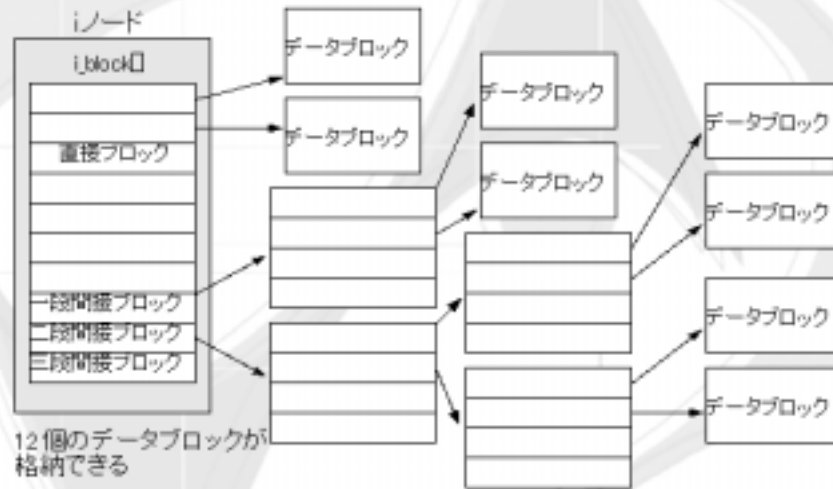
論理的なファイルと物理的なブロックとの対応関係を管理

- ext2 ファイルシステム
 - 古典的UNIXファイルシステムとほとんど同じ
- ext3
- reiserfs
- xfs
- jfs
- vfat
- iso9660

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



iノードとデータブロック



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



ブロックグループとスーパーブロック



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



iノードが保持している情報

i_mode: ファイル種別(通常ファイル, ディレクトリ, シンボリックリンク, キャラクタ型デバイス, ブロック型デバイス, パイプ, ソケット), ファイルモード

i_size: ファイルサイズ(バイト)

i_nlink: ファイルのリンク数(いくつかのディレクトリから参照されているか, ハードリンクされたファイルは2以上になる)

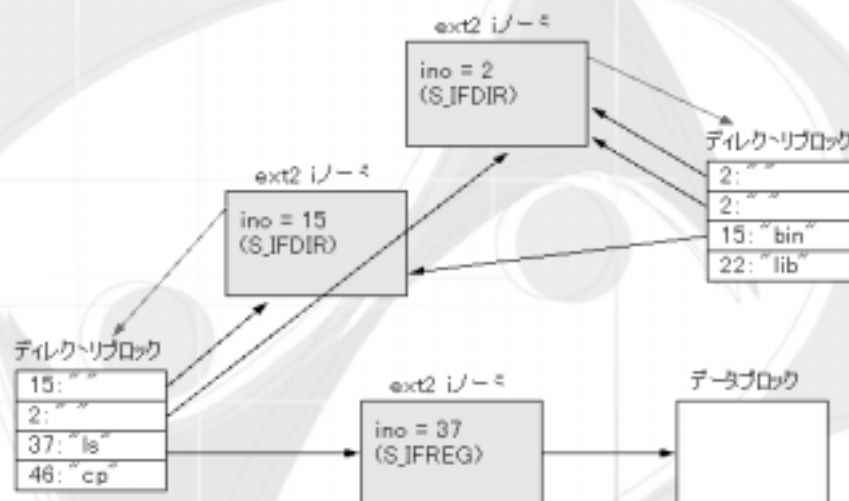
i_uid, i_gid: 所有者情報(ユーザID, グループID)

i_atime, i_ctime, i_mtime: ファイルのタイムスタンプ(アクセス時刻, iノード変更時刻, 更新時刻)

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



ディレクトリ階層構造



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



パス探索とキャッシュ

ファイル进行操作するには、パス名をiノードに変換する必要がある。

パス探索を高速化するためのキャッシュ

- ・iノードキャッシュ
- ・ディレクトリエントリキャッシュ

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



I/Oデバイスとファイル

- ・デバイスはデバイスファイルという形でVFSに用意される
 - ・キャラクタデバイス
 - ・ブロックデバイス
- ・メジャー番号・マイナー番号でデバイスを識別)
 - ・デバイスファイル名は(カーネルにとっては)関係ない
 - ・デバイスドライバ初期化時にメジャー番号との対応を登録

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.

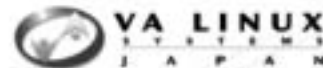


空間管理

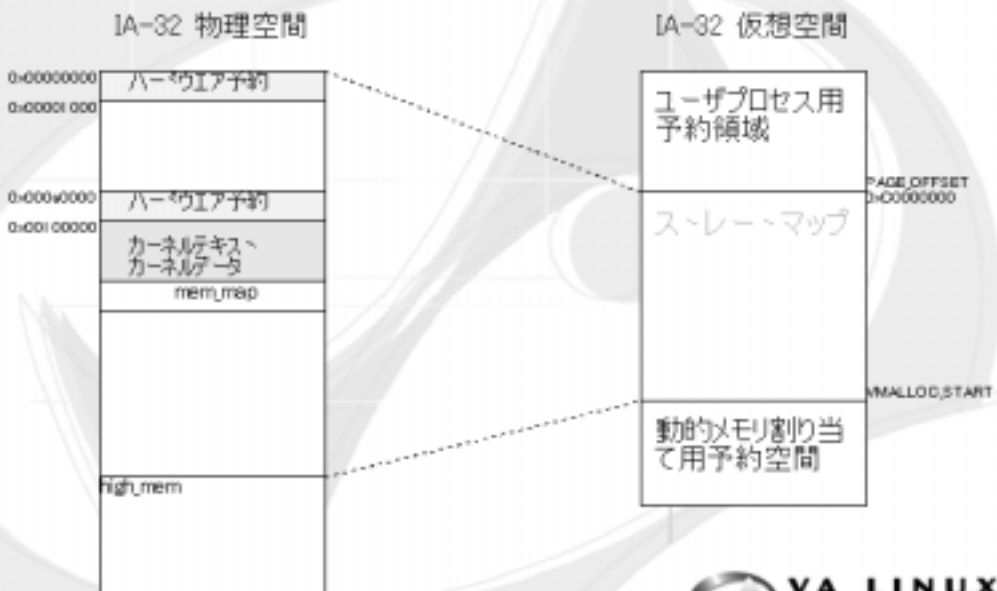
- 仮想メモリ空間(VM)を利用する
- ページ単位で管理する(4KB)
- 物理的なメモリと仮想メモリとの対応(ページテーブル)を利用して論物変換を行う(CPU)
- 実際の物理メモリよりも大きなサイズのプログラム群を動かすこともできる

| | |
|----------------------|----------|
| その他のメモリ管理コンポーネント | |
| スラバアロケータ | 動的空間割り当て |
| ページアロケータ ページ単位の管理 | |
| ハードウェア | |

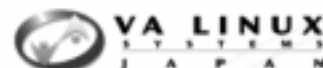
Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



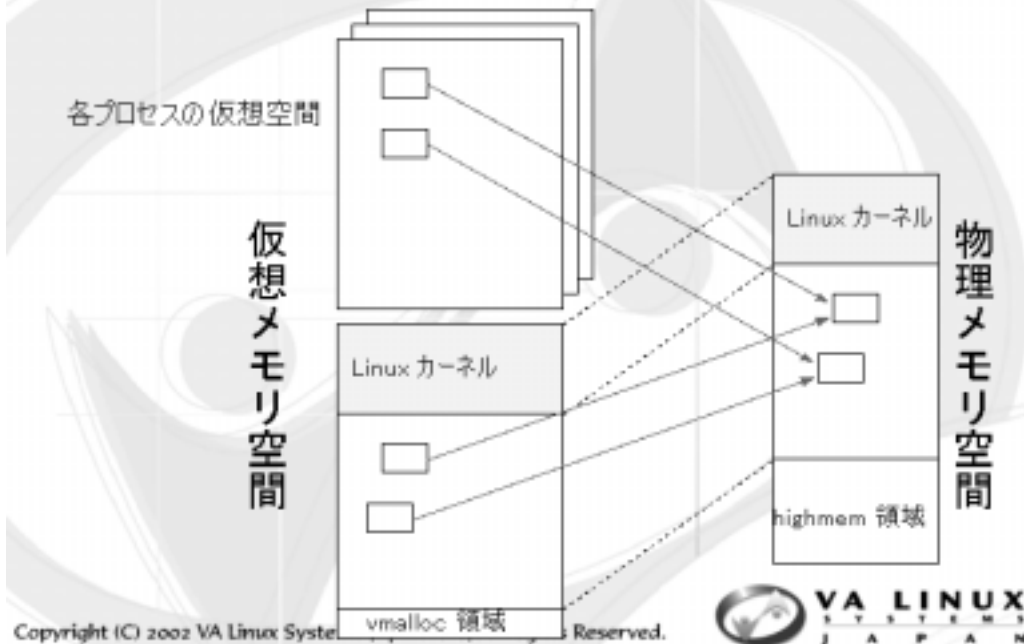
物理空間と仮想空間



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



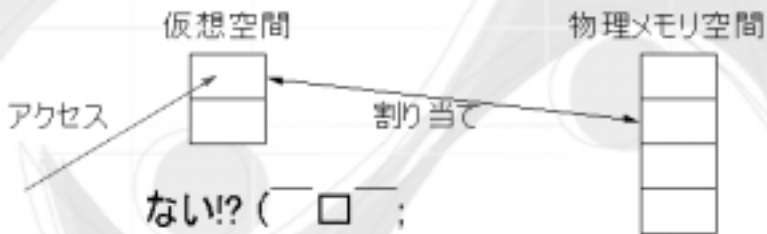
プロセスメモリ空間



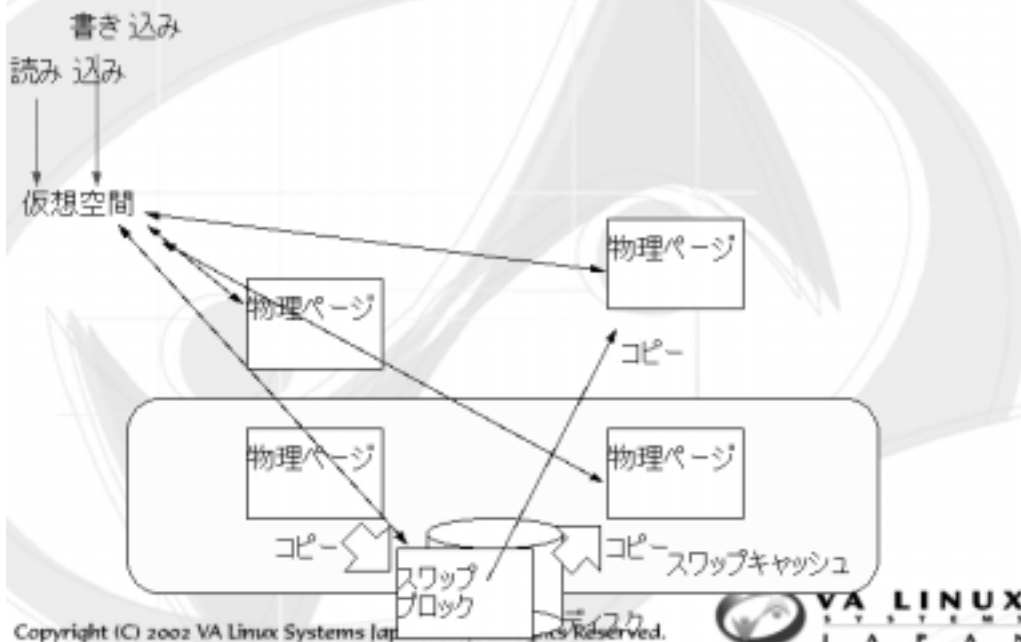
割り当て(デマンドロード)

できるだけ遅らせる

- ・デマンドロード
- ・コピーオンライト(あとで説明)



ページアウト・ページイン

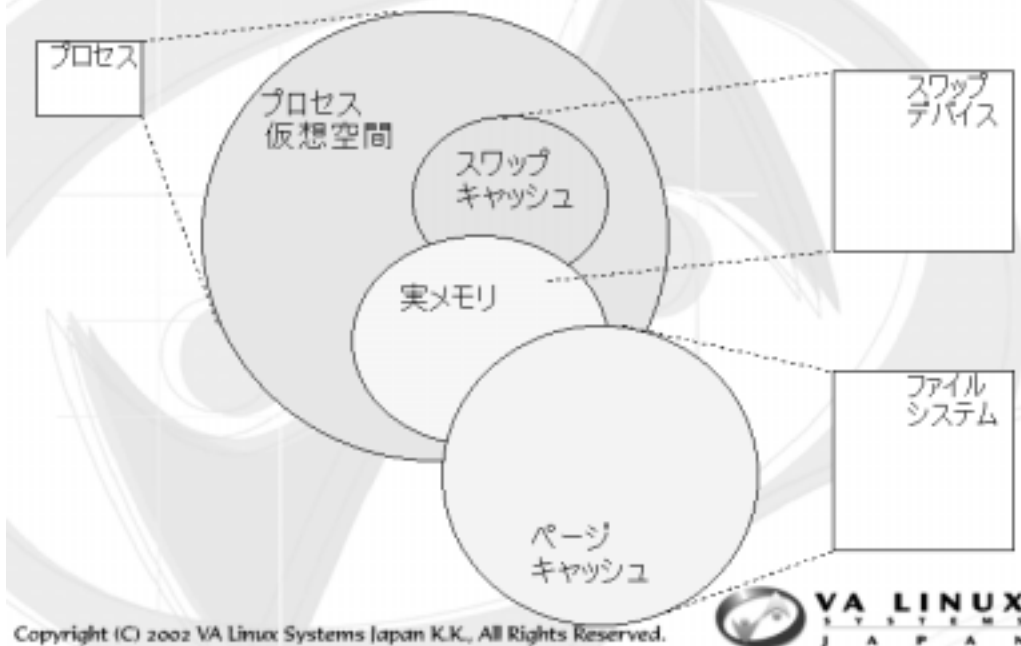


ページ解放について

- ・可能な限り空きメモリにせず、キャッシュとして利用
- ・利用価値の高いメモリは解放しない
- ・高負荷時には早めに多めに確保

カーネルスレッド kswapd による監視

キャッシュとの関係

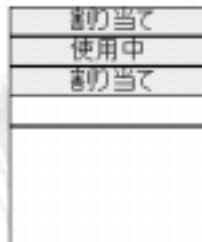


ページアロケータ

- ページの管理
- ページごとの管理情報
 - 参照カウント
 - 仮想マップ先のアドレス
 - アクティブか?
 - ...

バディシステム

- ページアロケータの管理方法
 - 無作為にページを割り当てると、メモリのフラグメンテーションが起きやすい
- 2のべき乗個単位でページを管理する



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.

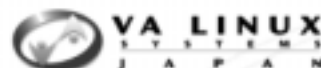


スラブアロケータ

- ページよりも小さい単位でのメモリ割り当て
- `/proc/slabinfo` 参照



Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



その他のトピック

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



SMPカーネル

- SMP…Symmetric Multi Processing
 - すべてのCPUがメモリを共有
 - すべてのメモリに均一なアクセスが可能
 - メモリバスの速度が性能に影響する

リソースの競合状態を引き起こさないように
特別な配慮が必要

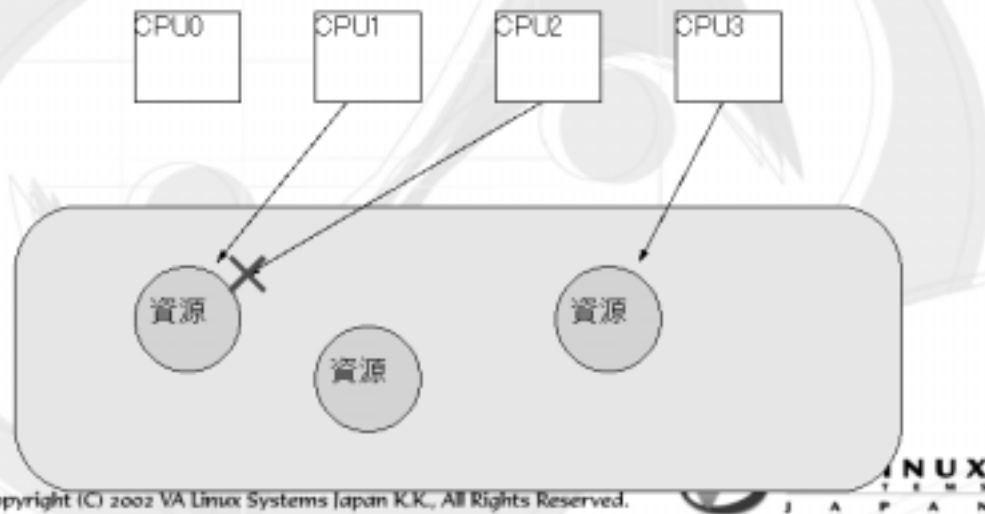


Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



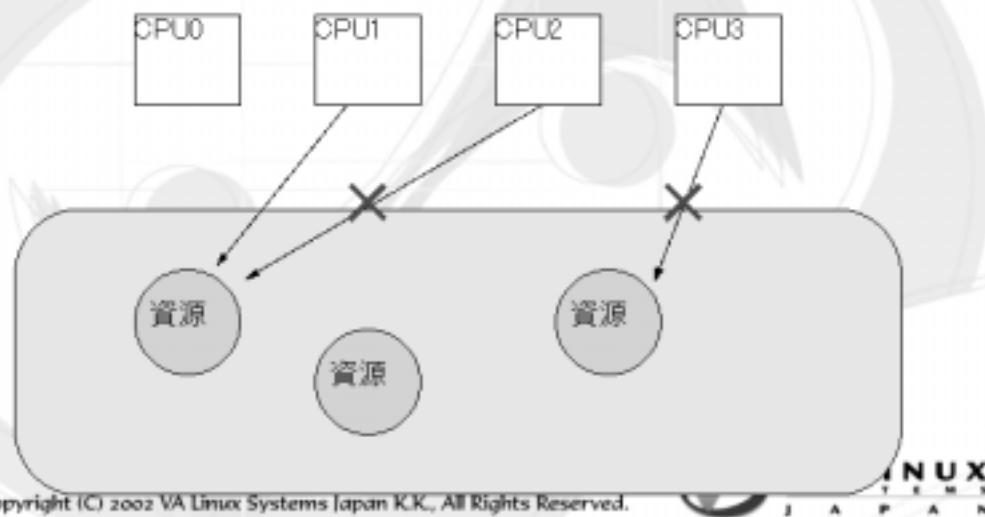
排他制御の仕組み 1

・スピンロック



排他制御の仕組み 2

・カーネルロック(ジャイアントロック)



モジュール

- 動作中のLinuxカーネルに対して、新しい機能を追加するしくみ
 - デバイスドライバ
 - ファイルシステム
- 自動組み込みが可能 (kmod)
- カーネルの機能とユーザコマンドで利用する
 - insmod rmmod modprobe など

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



モジュールの組み込み・削除

- 組み込み
 - モジュール配置領域の予約
 - vmallocで動的に領域確保
 - 未解決の参照の解決
 - モジュールの組み込み
 - init_module()システムコール
- 削除
 - delete_module()システムコール
 - 依存されていなければ削除してvfreeでメモリを解放

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



ソースツリーについて

| | |
|----------------|------------------------|
| /usr/src/linux | |
| 6MB | Documentation (ドキュメント) |
| 26MB | arch (アーキテクチャ依存コード) |
| 73MB | drivers (ドライバ) |
| 10MB | fs (ファイルシステム) |
| 20MB | includes (インクルードファイル) |
| 0.4MB | kernel (カーネルコアコード) |
| 0.4MB | mm (メモリ管理) |
| 6MB | net (ネットワーク) |
| 0.4MB | scripts (スクリプトファイル) |

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



参考文献・資料

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



参考書籍・記事

- 詳解Linuxカーネル(オライリージャパン, 2001)
Daniel P. Bovet, Marco Cesati 著
高橋浩和, 早川仁 監訳
- Linux Japan連載 Linuxカーネル2.4の設計と実装
→LJ最終号にPDFで収録

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.



参考URLなど

- <http://www.kernel.org/>
Linuxカーネルの一次配布もと
- Linux Kernel Documents
カーネルソースツリーのDocumentsディレクトリ
JF(<http://www.linux.or.jp/JF/>)に日本語訳がある

Copyright (C) 2002 VA Linux Systems Japan K.K., All Rights Reserved.

