

組込みエンジニアのためのLinux入門 仮想メモリ編

2007.2.22 original

2012.10.27 revival

小林哲之

注釈

- このスライドは2007年に作成して、CELFテクニカルジャンボリーで発表したものです。
- <http://tree.celinuxforum.org/CelfPubWiki/JapanTechnicalJamboree13>
- 内容的には現在でもほぼ通用するので再演します。

このスライドの対象とする方

- 今までずっと組込み機器のプロジェクトに携わってきて最近OSにLinuxを使っている方々

このスライドの目的

- Linuxの仮想メモリの仕組みを理解し現在のプロジェクトに役立てる。
 - 仕組みを知らなくてもプログラムは動くが、性能を引き出すためには仕組みの理解が重要。

まず基本の概念から

- 仮想～、論理～ virtual, logical
 - 仮想アドレス、論理デバイス、論理セクタ、仮想マシン
 - あたかも ... のように扱う
- 実～、物理～ real, physical
 - 実アドレス、物理デバイス、物理セクタ
 - そのもの、そのまんま

仮想化：あたかも・・・ 実は

- あたかも巨大のようだが、実は少ない。
- あたかも平らのようだが、実は凸凹。
- あたかもたくさんのようにだが、実はひとつ。
- あたかも占有しているようだが、実は共有。

仮想化は複雑さや個々に依存することを
隠蔽するマジック。

マジックなので種も仕掛けもある。

＝ 実と仮想の対応付け（マッピング）

あたかもそう見えるように変換している。

仮想化の代償

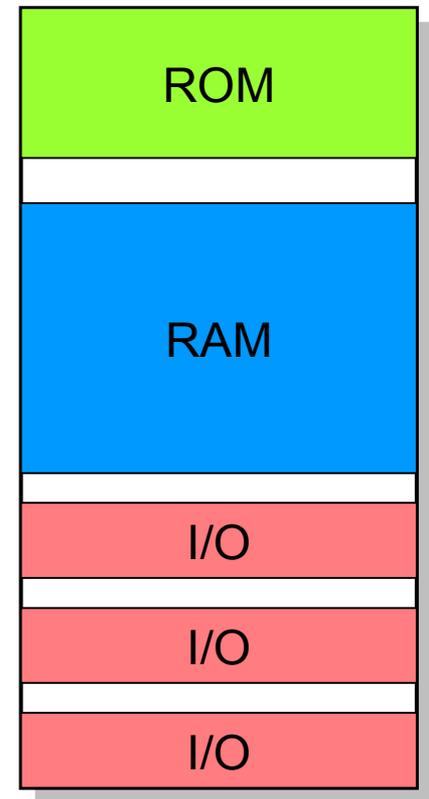
- マイナス面よりもプラス面が勝っているから仮想化を行うわけだが、常にプラスとは限らない。

物理メモリと仮想メモリ

- いまままでの大抵の組込み機器プロジェクトでは物理メモリしか扱うことがなかった。
- 最近では組み込みシステムの規模の増大化に伴ってLinuxやWindow CEなどPC向けOSの流れを持つOSを使用することが多くなってきた。これらのOSは仮想メモリシステムを備えている。
- このスライドではLinuxについて説明する。

物理メモリ

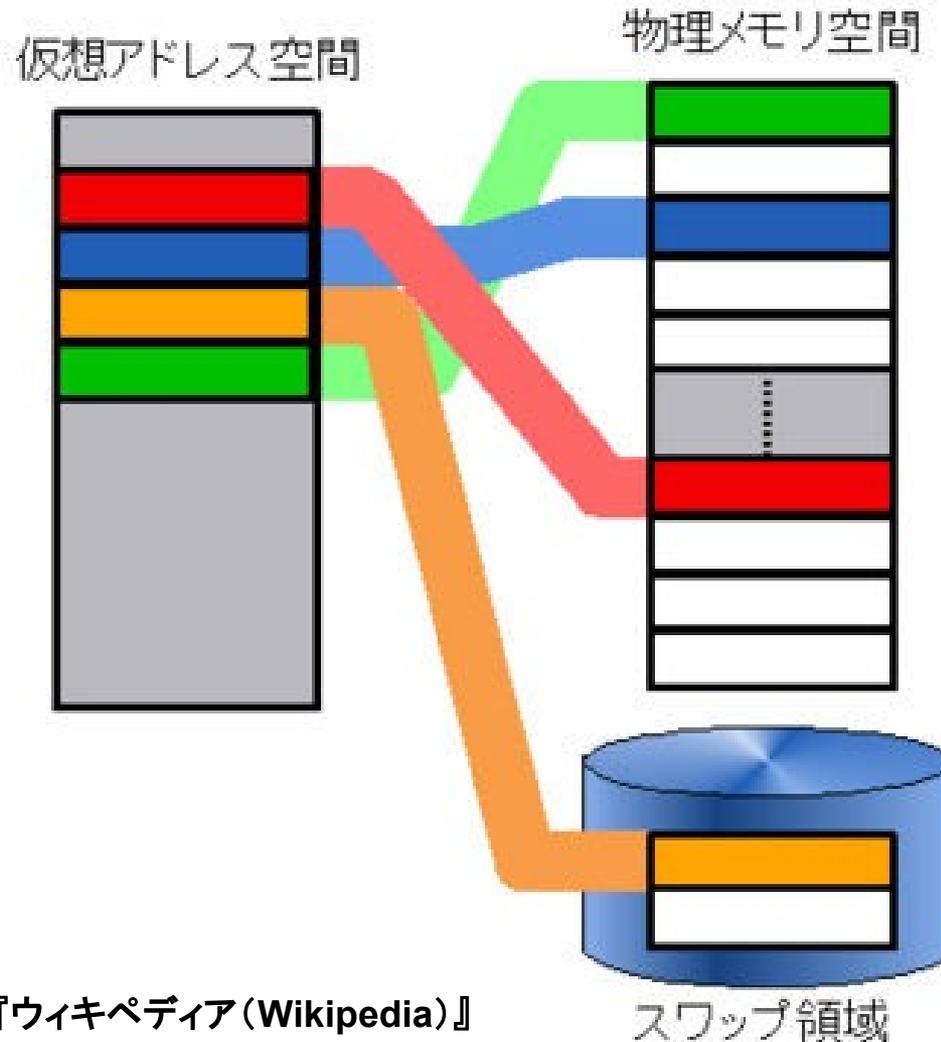
- 単一のメモリ空間
- 機器ごとにROM, RAM, I/Oの実装アドレスが異なるのでそれを意識してプログラムする



仮想メモリ

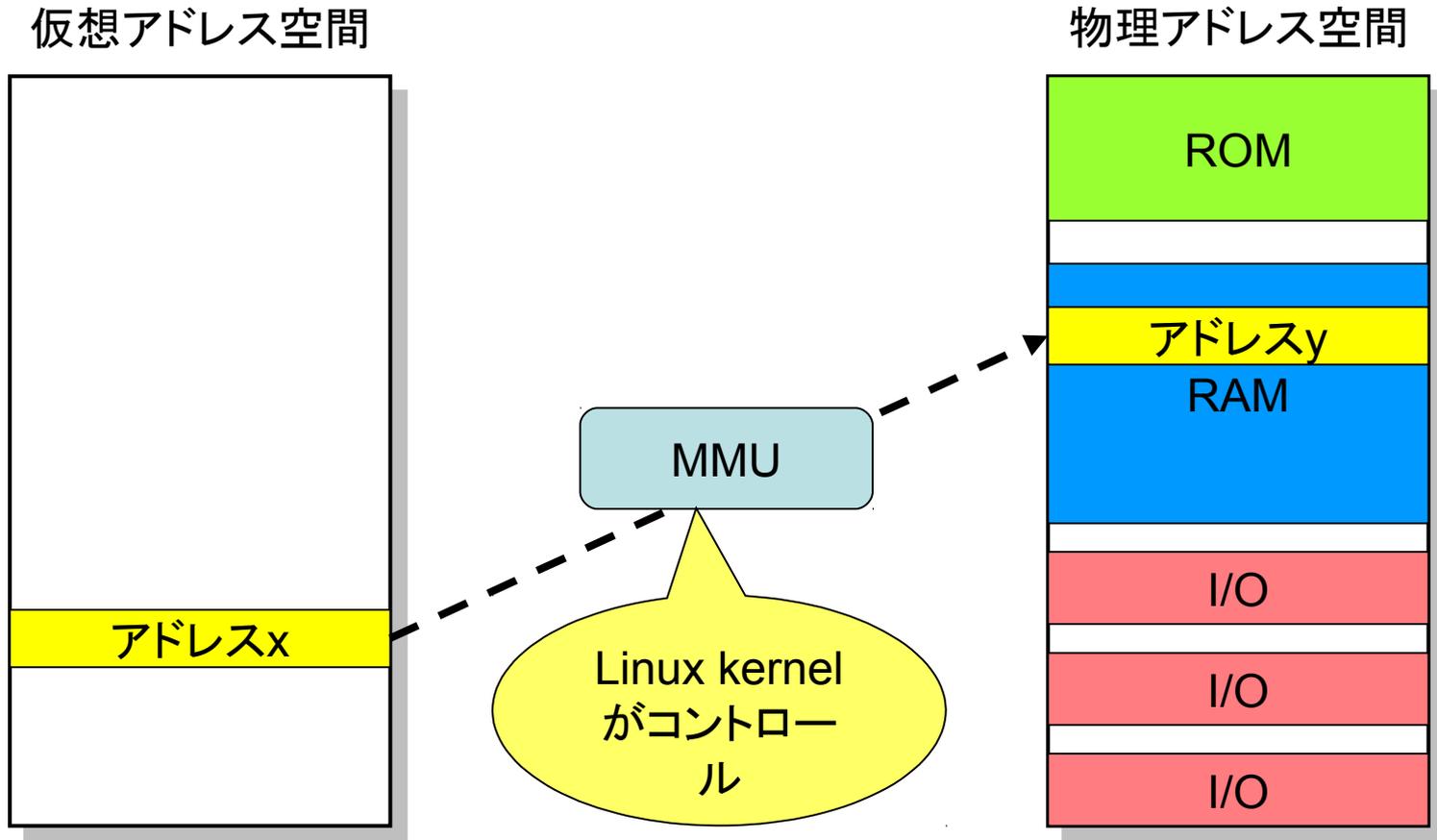
- 利点
 - ユーザープログラムは実際のメモリマップ(実装アドレス、実装サイズ)に依存しなくなる。
 - 不連続な物理メモリの断片を連続する仮想メモリとして利用できる。
 - メモリ保護: バグによって無関係の部分のメモリが破壊されることを防止できる。
- 新しい概念の導入
 - アドレス変換
 - 多重メモリ空間
 - デマンドページング

仮想メモリの概念図

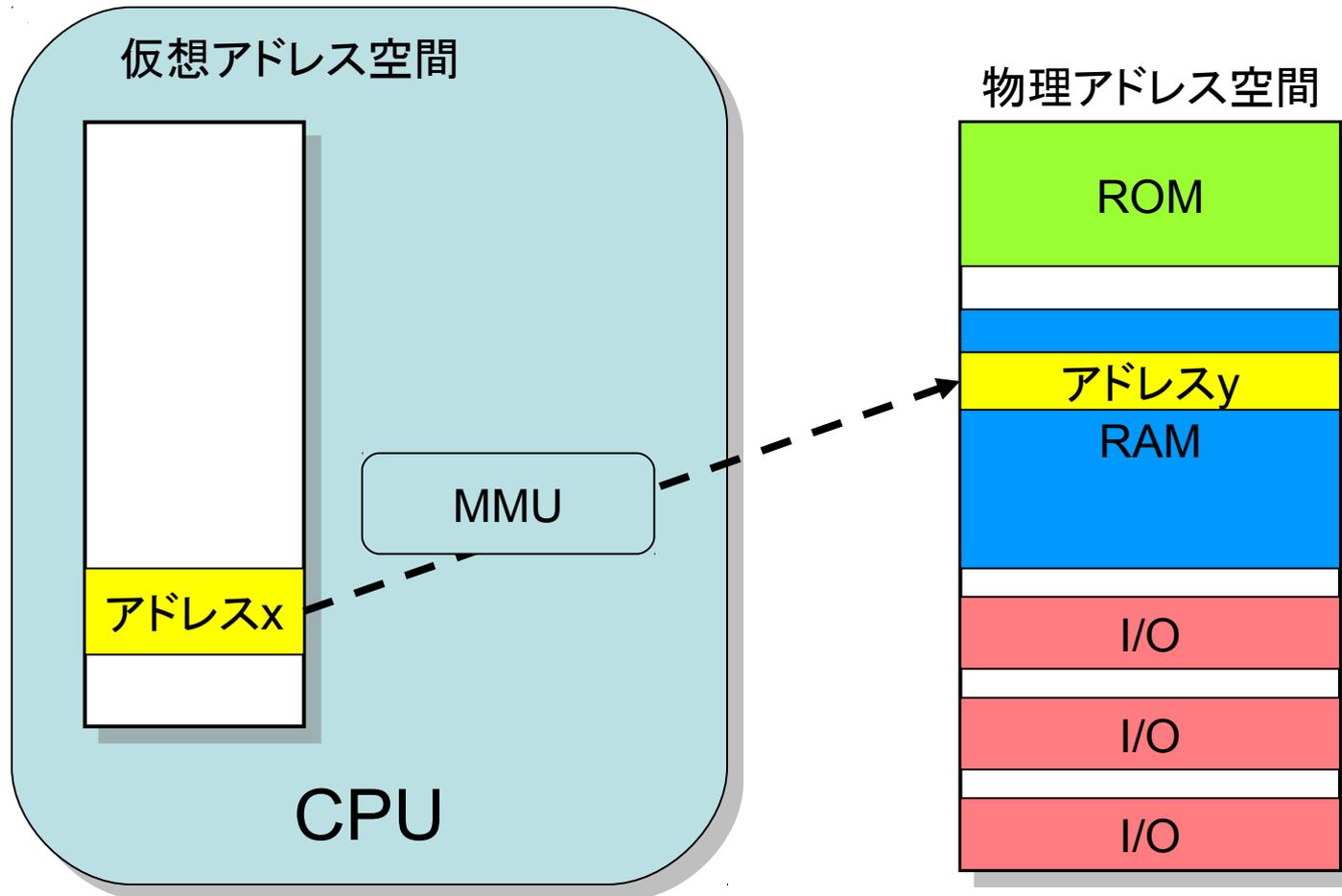


出典: フリー百科事典『ウィキペディア (Wikipedia)』

アドレス変換

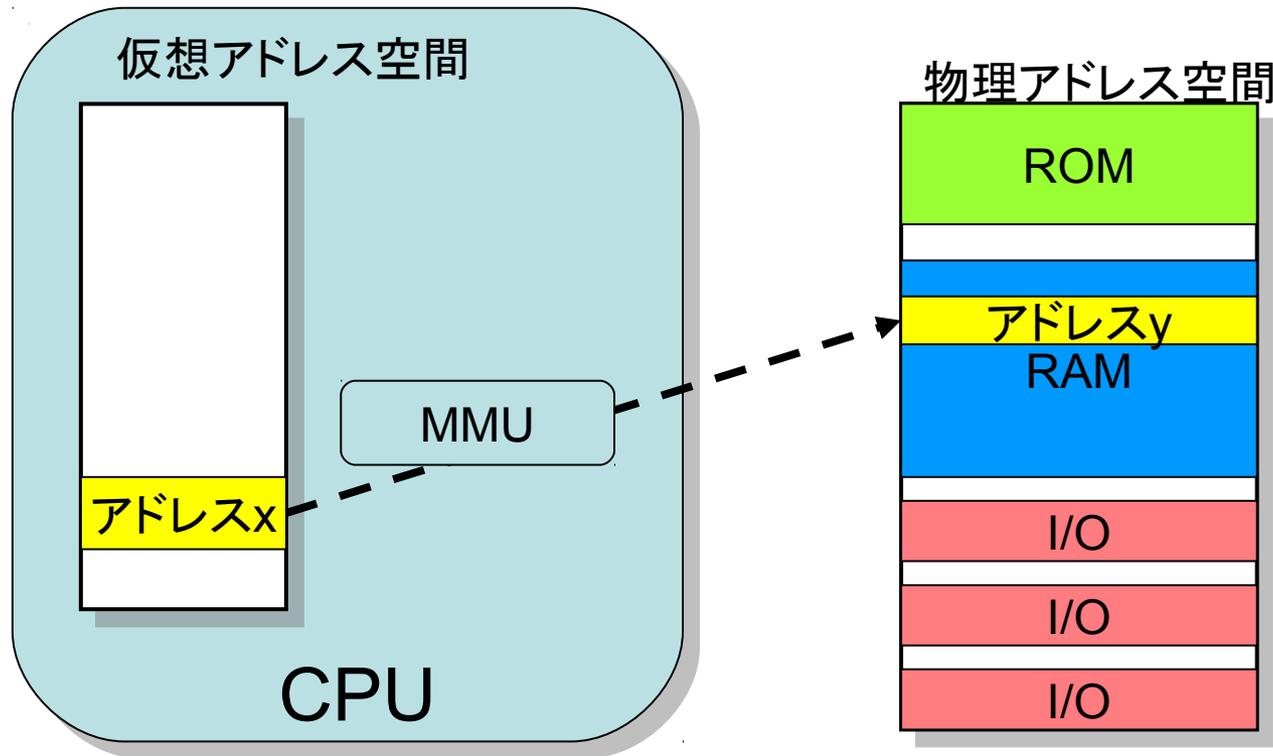


仮想メモリはCPUの中だけ



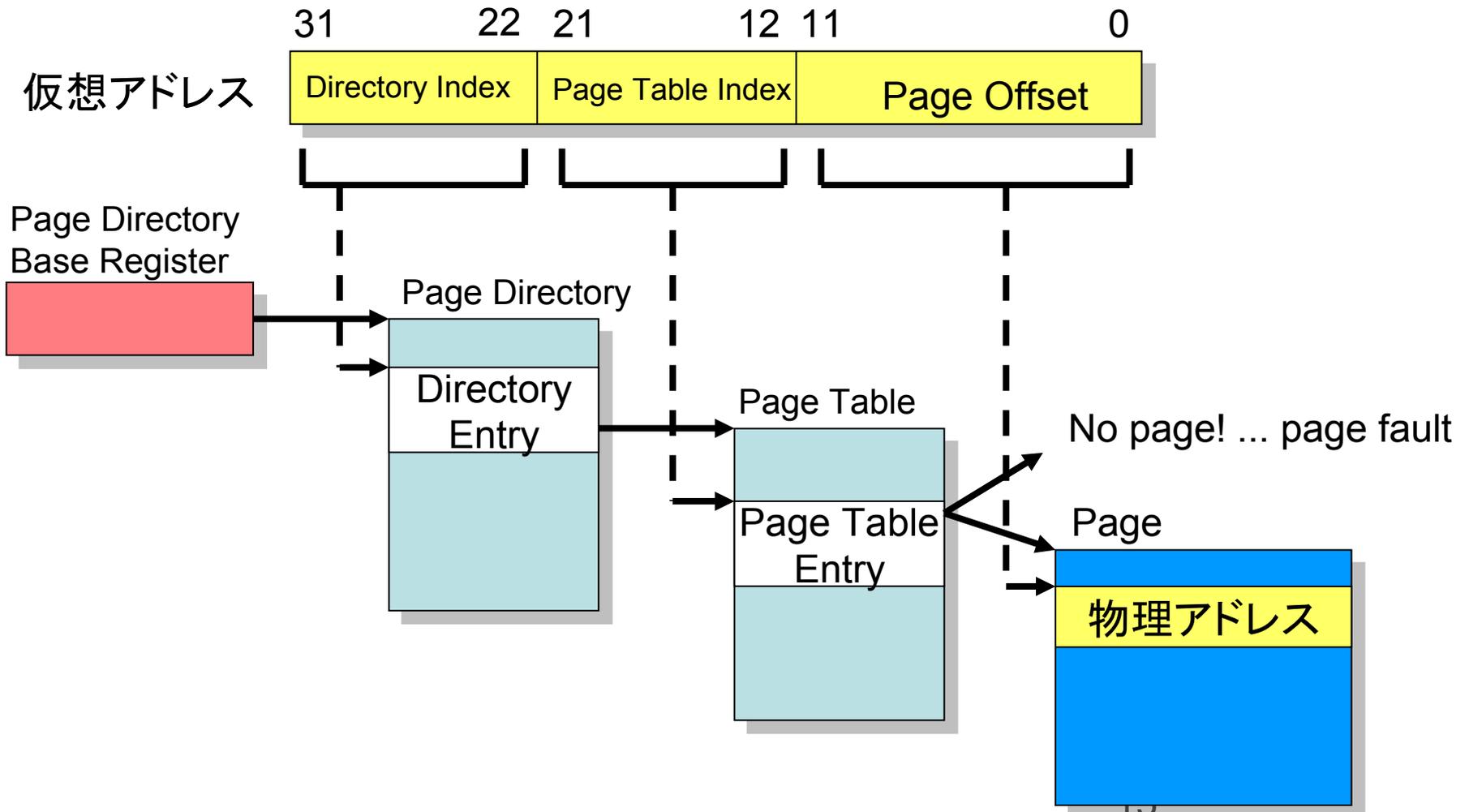
CPUから外にでてくるアドレスバスにのるのは物理アドレスだけ。
ロジアナでは仮想アドレスは観測できない。

ユーザープログラムは仮想アドレスだけ



物理アドレスを扱うのはカーネルモードだけ。
つまりカーネル本体とデバイスドライバ。

MMUでのアドレス変換



TLB

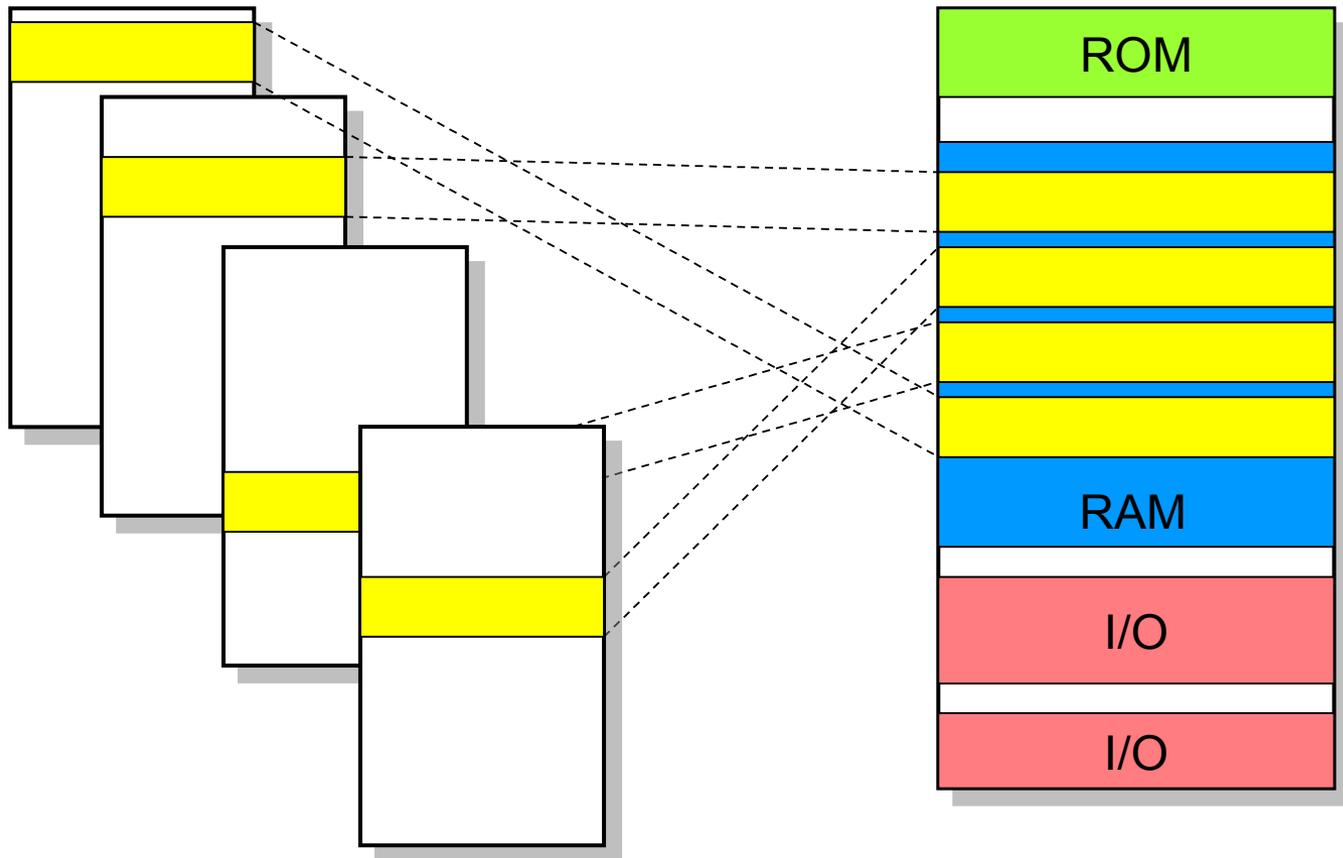
仮想アドレスページ	物理アドレスページ
	...

- Translation Lookaside Buffers
- 仮想アドレスをkeyとして物理アドレスを得るハッシュテーブルのようなもの
- 大抵のアドレス変換はTLBにヒットするので実際にPage Directory, Page Tableをアクセスせずに済む。

多重メモリ空間

プロセスごとに独立した仮想メモリ空間

物理アドレス空間



デマンドページング

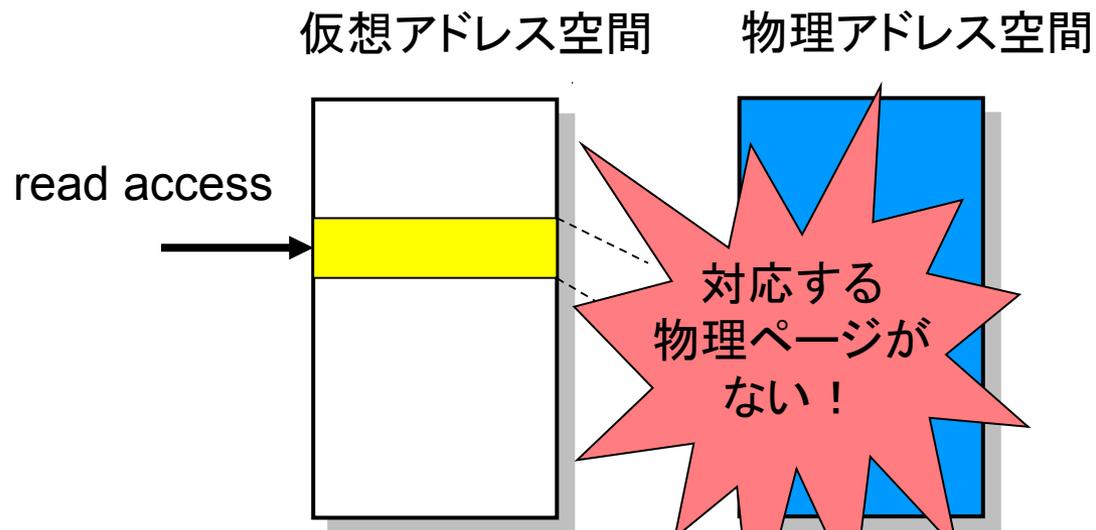
- ページ単位でマッピングされる
 - ページのサイズはたいていは4Kbytes
- 2段階で行われる
 1. 仮想メモリの割り当て(mmap)
管理台帳に登録するだけ
 2. 実際にアクセスがあったときに初めて
そのページに物理メモリが割り当てられる

アクセスのないページには物理メモリが割り当てられないので

仮想メモリサイズ \geq 実際に必要な物理メモリサイズ

デマンドページングの動作例

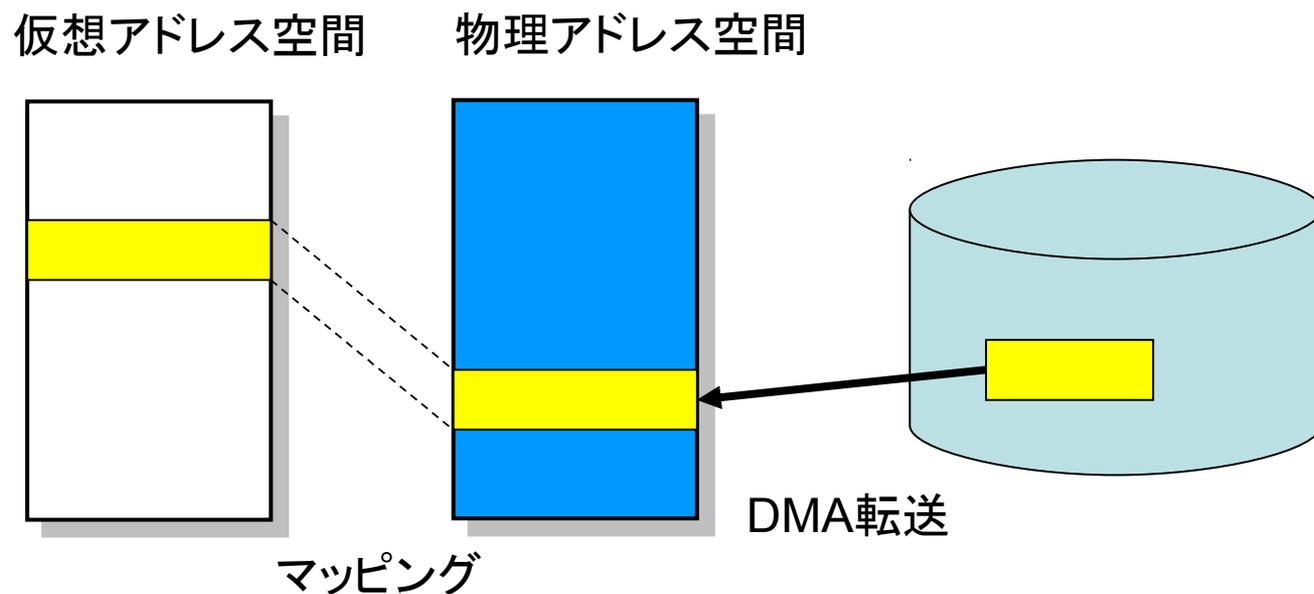
(1)



ページフォールト発生
カーネルモードへ

デマンドページングの動作例(続)

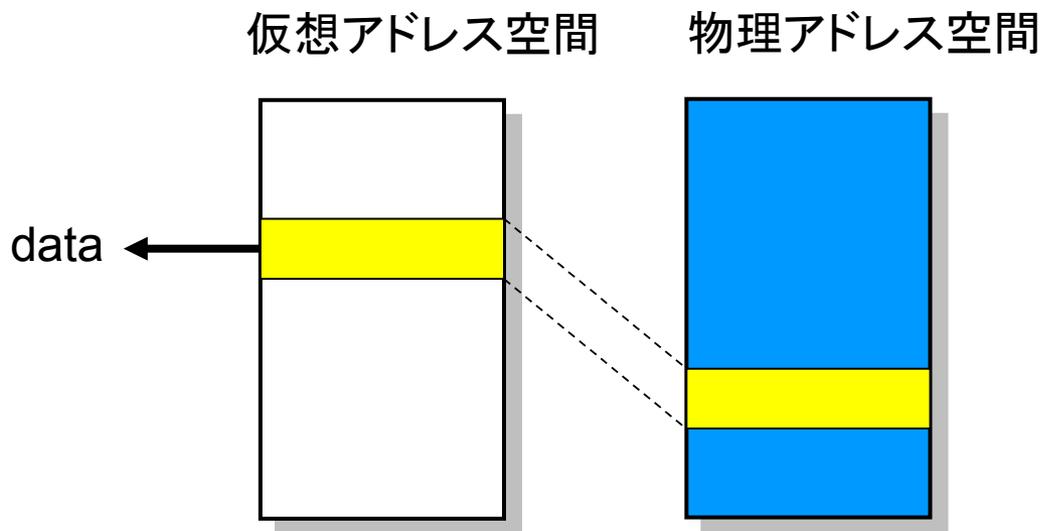
(2)



カーネルがデータをロードして物理アドレスをマッピングする。

デマンドページングの動作例(続)

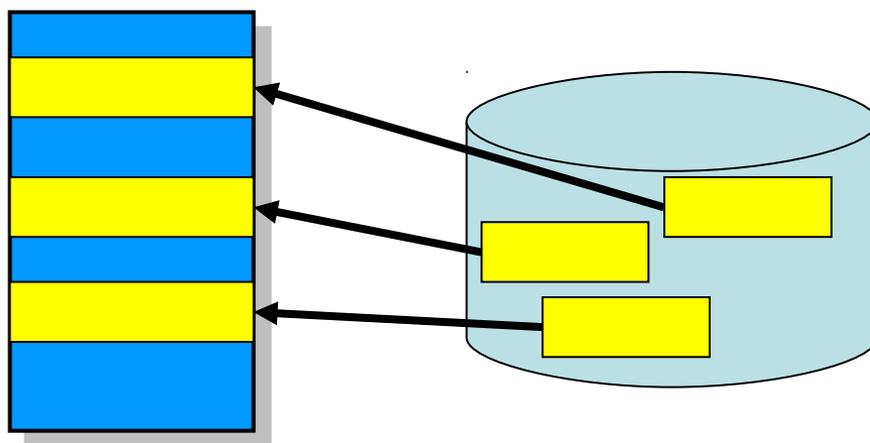
(3)



ユーザーモードに復帰。
ユーザープログラムからは何事もなかったように
データが読める。... でも実際に時間はかかっている。

ページキャッシュ

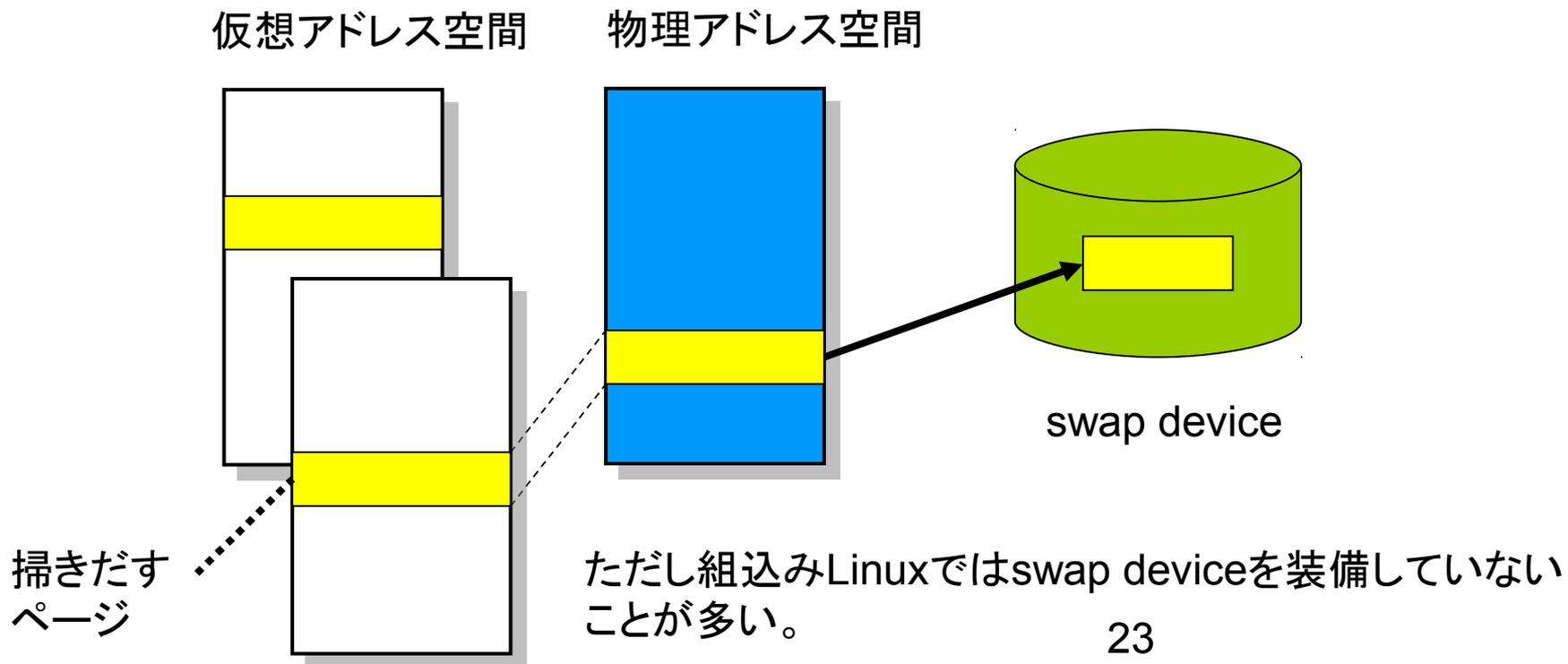
物理アドレス空間



ディスクからの読んだ内容はメモリに余裕がある限り保持しておく。
シーケンシャルにアクセスされる場合が多いので数ページ分を
まとめて先読みを行う。
そのため(2)では毎回ディスクアクセスが発生するとは限らない。

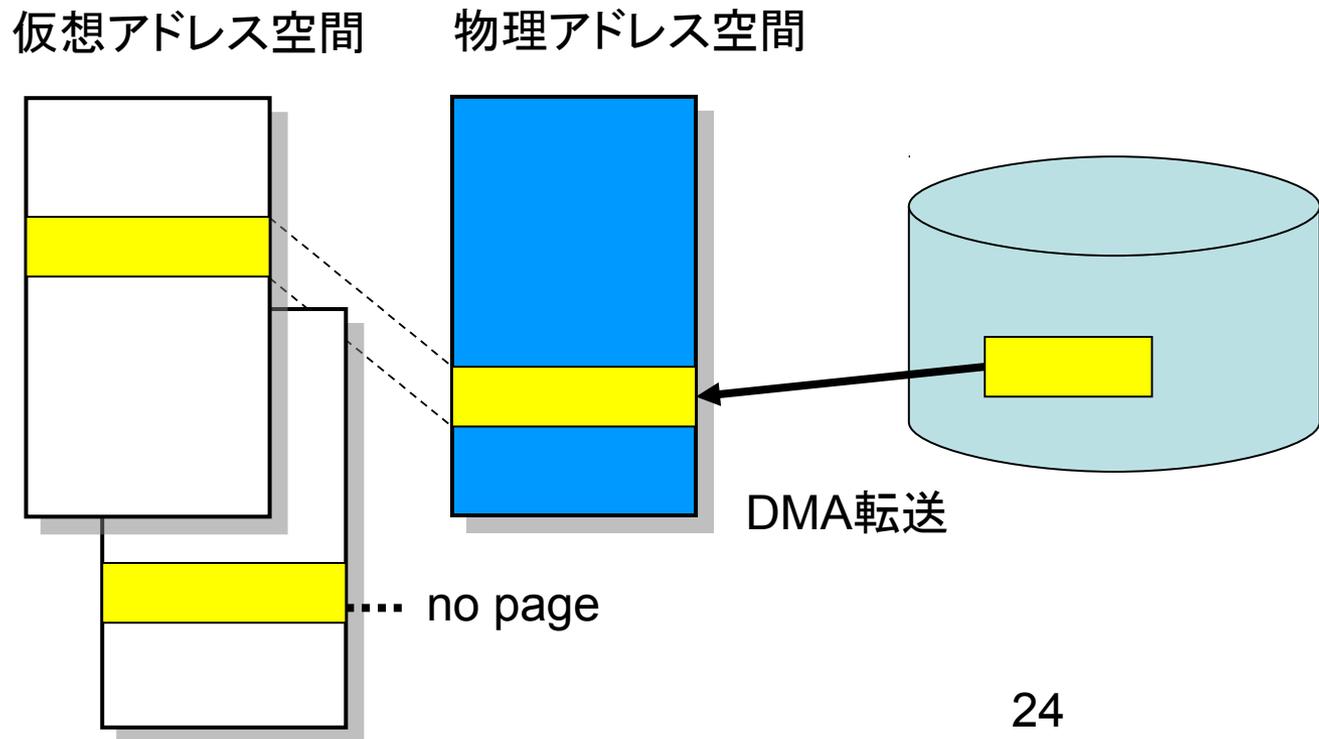
ページアウト

(2)で物理メモリの空きがなかった場合、使用頻度の低いと思われるページを解放する。そのページの内容が変更されていなければそのまま破棄。変更されていればスワップデバイスに掃きだす。

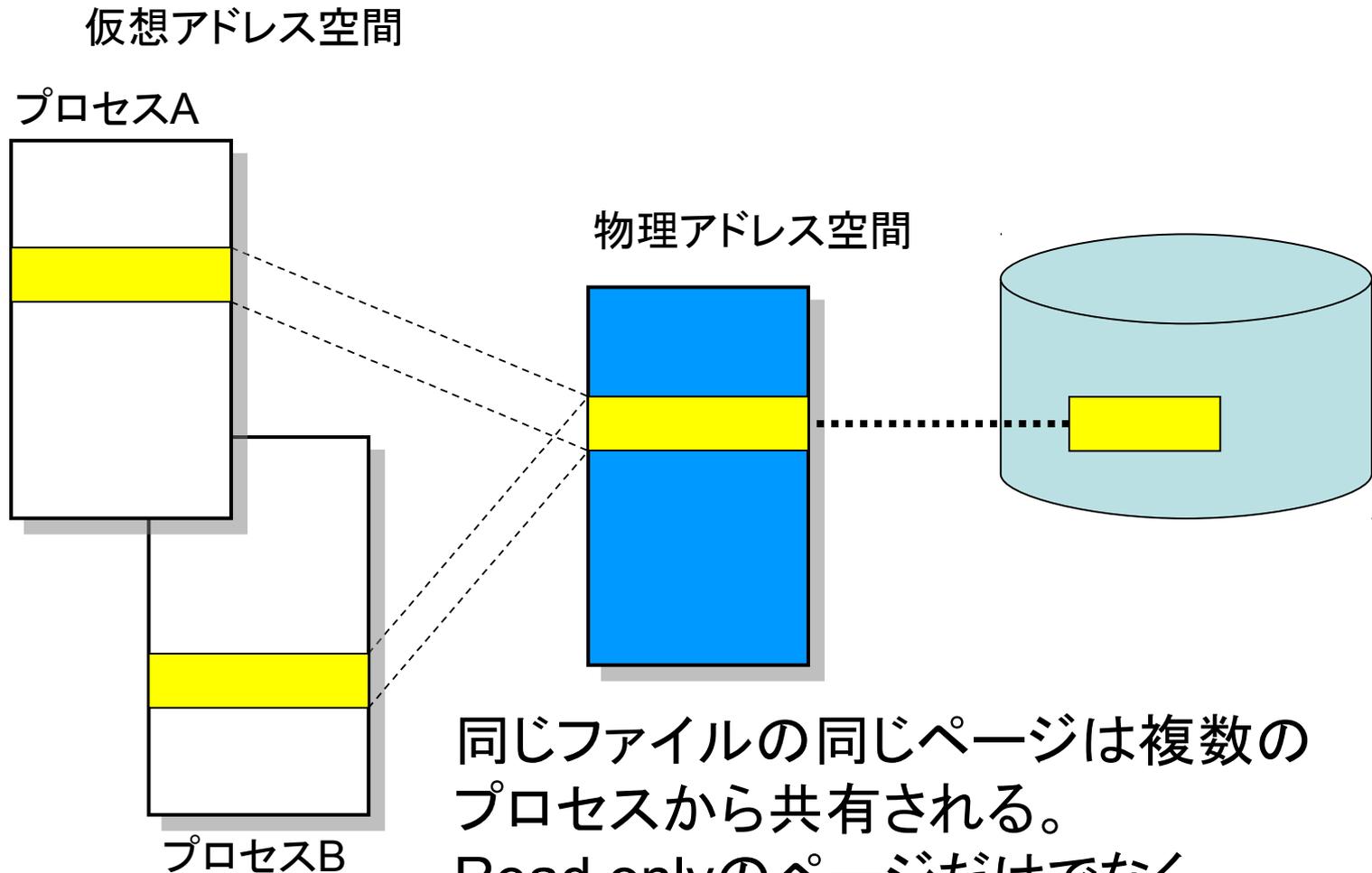


ページアウト(続)

解放した分のページを使って要求されたページの割り当てを行う。
このような「お手玉」をすることで実際に搭載されている物理メモリサイズよりも大きなサイズの仮想メモリを扱うことができる。



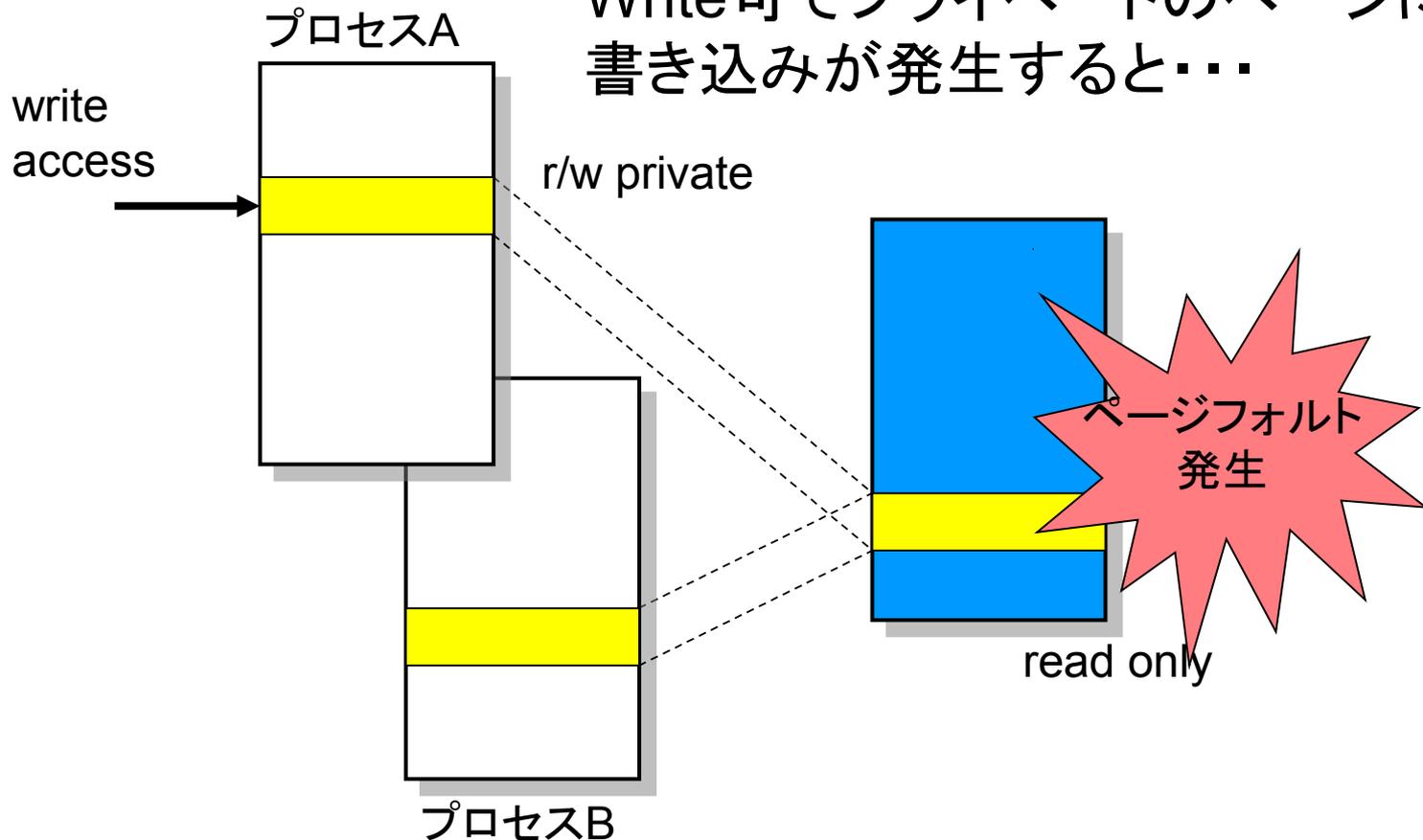
ページの共有



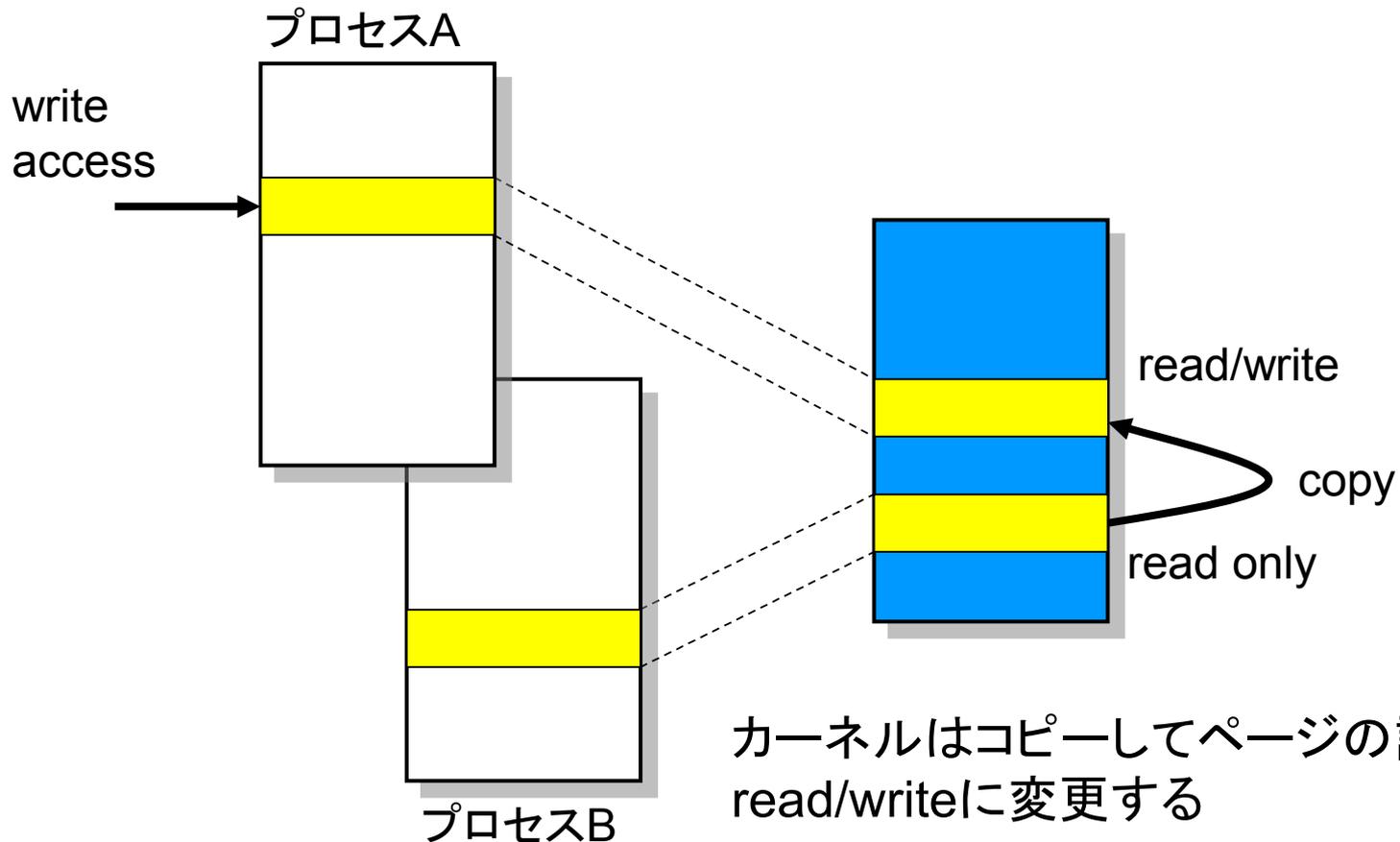
同じファイルの同じページは複数のプロセスから共有される。
Read onlyのページだけでなく
write可のページも。

コピーオンライト

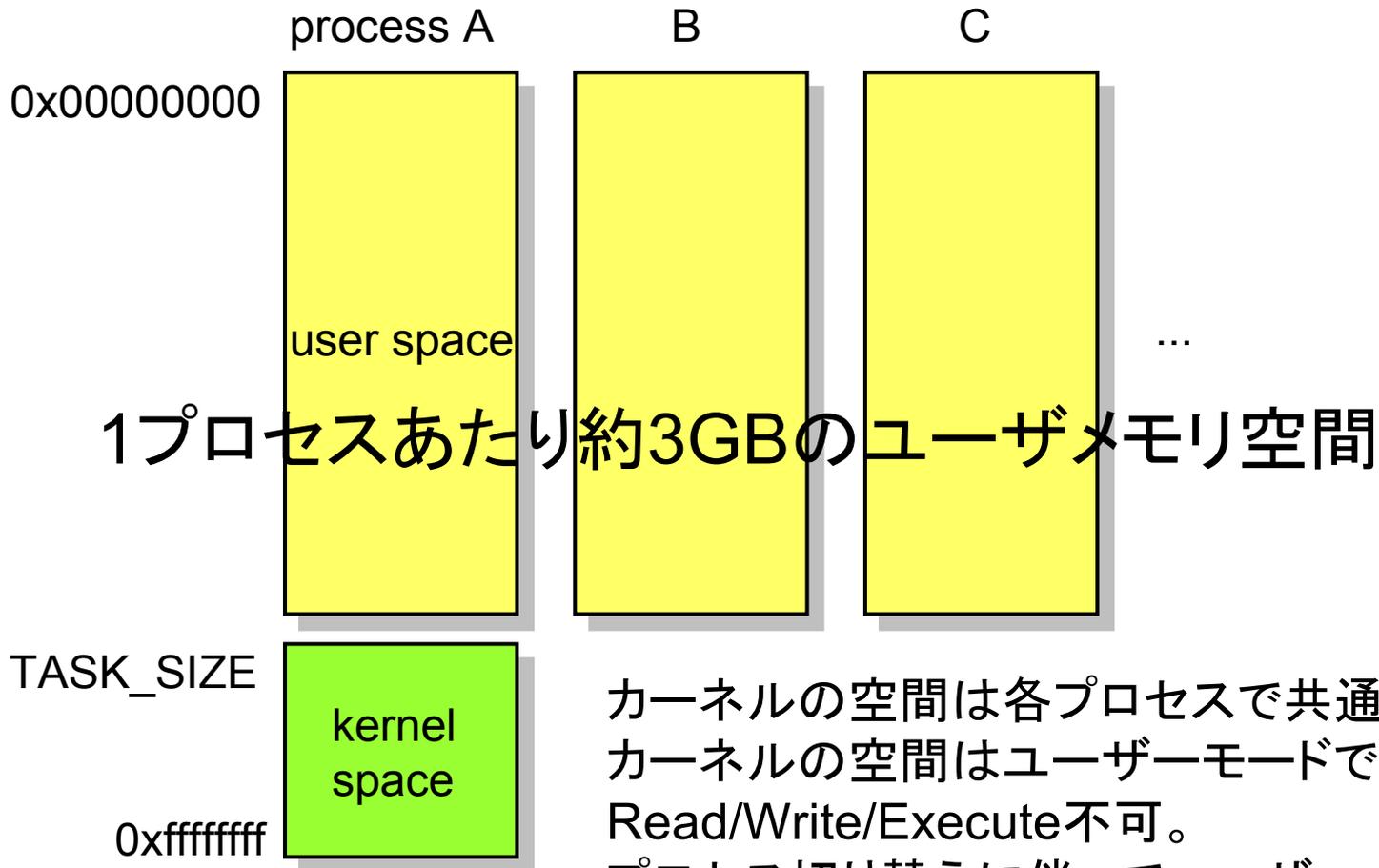
Write可でプライベートのページに書き込みが発生すると...



コピーオンライト(続)



プロセスのメモリ空間



TASK_SIZEはi386では0xc0000000
ARMでは 0xbf000000

カーネルの空間は各プロセスで共通。
カーネルの空間はユーザーモードでは
Read/Write/Execute不可。
プロセス切り替えに伴ってユーザーメモリ
空間が切り替わる。

ユーザプロセスのメモリ空間の実例 (詳細)

```
cat /proc/<PROCESS_ID>/smaps
```

```
.....  
0011e000-0024a000 r-xp 00000000 fd:00 15172740 /lib/libc-2.4.so  
Size:                1200 kB  
Rss:                 136 kB  
Shared_Clean:       136 kB  
Shared_Dirty:        0 kB  
Private_Clean:      0 kB  
Private_Dirty:      0 kB  
0024a000-0024d000 r-xp 0012b000 fd:00 15172740 /lib/libc-2.4.so  
Size:                12 kB  
Rss:                 8 kB  
Shared_Clean:       0 kB  
Shared_Dirty:       0 kB  
Private_Clean:     0 kB  
Private_Dirty:     8 kB  
0024d000-0024e000 rwxp 0012e000 fd:00 15172740 /lib/libc-2.4.so  
Size:                4 kB  
Rss:                 4 kB  
Shared_Clean:       0 kB  
Shared_Dirty:       0 kB  
Private_Clean:     0 kB  
Private_Dirty:     4 kB  
.....
```

RSS = 物理メモリサイズ

システムコールmmap

```
#include <sys/mman.h>

void *mmap(void *start, size_t length, int prot, int flags,
           int fd, off_t offset);

int munmap(void *start, size_t length);
```

- ファイルやデバイスをメモリにマップ/アンマップする
- 引数 *prot*
 - PROT_NONE または PROT_EXEC, PROT_READ, PROT_WRITEのOR演算
- 引数 *flags*
 - MAP_FIXED, MAP_SHARED, MAP_PRIVATE, MAP_ANONYMOUS, ...

mmapのtips

- MAP_FIXEDを指定しなければカーネルが空いているページをさがしてくれる。
- MAP_FIXEDを指定したときに既存のページと重なっていたら、そのページは内部的にmunmapされる。
 - なのでこのオプションは通常は使用しない。
- ファイルのオフセットはページサイズの整数倍でなければならない。
- mmapとmunmapのアドレス、サイズは一致していなくてもよい。

mmapの使い方(1)

- 巨大サイズのmallocの代用
 - コンパクションなどデータのコピーが発生しない。
 - malloc/freeと違ってmunmapするときのaddr, sizeはmmapで確保したときと異なってもよい。
 - まとめて1回のmmapで確保して少しずつ分割してmunmapで返却するのもあり。
 - glibcのmallocの実装ではある一定以上のサイズのmallocはmmapを呼び出す。
 - `DEFAULT_MMAP_THRESHOLD = (128*1024)`

mmapの使い方(2)

- 高速なファイルアクセス
 - read, writeのシステムコールでは内部で物理ページにバッファリングしている。そこからユーザの指定した配列にコピーしている。
 - mmapを使うことで直接ページにアクセスできるようになるのでデータのコピーを減らすことができる。
 - Java 1.4の `java.nio.MappedByteBuffer`

mmapの使い方(3)

- プロセス間の共有メモリ
 - 複数のプロセスから同じファイルをR/W可sharedでマッピングする。
 - IPCの共有メモリのシステムコール (shmget, shmat,..)は内部で同様のことを行っている。

mmapの使い方(4)

- 物理メモリ、I/Oポートのアクセス
 - デバイスファイル `/dev/mem` をマッピングすることでユーザーモードで物理メモリ空間をread/writeすることが可能。
 - `/dev/mem`をアクセスするにはrootの権限が必要。

まとめ

- 仮想メモリの使用量と物理メモリの使用量は異なる。実際に問題になるのは物理メモリの使用量。
- 仮想メモリのオーバーヘッドはいつ発生するのかを意識する。
 - TLBミス
 - ページフォルト
- システムコールmmapの活用。

参考文献

- Linux kernel ソース
<http://www.kernel.org/>
- GNU C ライブラリのソース
<http://www.gnu.org/software/libc/>
- “**詳解Linuxカーネル 第2版**”
オライリージャパン
- “Linuxカーネル2.6解説室”
SoftBank Creative
- Linux manコマンド
- その他たくさんのWEB検索結果

2007年当時

おまけ：~~最近~~の話題

- CELFのBootTimeResourcesより
 - KernelXIP
 - ApplicationXIP
 - (DataReadInPlace)

NORフラッシュが使われなくなったためあまり聞かなくなった。
- CELFのMemoryManagementResourcesより
 - Huge/large/superpages
 - Page cache compression

→ zramの名前でメインライン化。
3.8ではstaging卒業の見込み。