# 6

# Virtual Storage

*Nil posse creari de nilo*
*– Lucretius, De Rerum Natura, I, 155*

## 6.1 Introduction

In this chapter, mechanisms to support virtual storage will be modelled. Virtual storage affords a considerable number of advantages to the operating system designer and user. Virtual storage is allocated in units of a page and pages can be collected into independent segments; virtual storage defines clear boundaries between address spaces so that each process can maintain its own address space. This clearly provides a measure of protection between processes but requires additional methods for exchanging data between processes.

In virtual storage systems, main store is shared between processes in the usual way but is defined as a sequence of *page frames*, each a block of store one page long. The storage belonging to each process is swapped into main store when required and copied to a paging disk (or some other mass-storage device) when not required. Strategies for selecting pages to copy to the paging disk and for determining which page to bring into main store must be defined.

## 6.2 Outline

The storage system to be designed is to have the following features:

- The virtual store should have four segments: one each for code, data, stack and heap.
- The system uses *demand paging* with reference counting for victim selection.
- Pages can be shared (and unshared) between processes.
- Segments can be shared (and unshared) between processes;

- Storage should be mapped in the sense that disk blocks can be directly mapped to main-store pages and vice versa.
- Message passing will be used for IPC.

The virtual storage system is composed of:

- A page-fault handler. This is invoked when a page fault occurs. It determines the identity of the *logical* page that caused the page fault. It invokes the page-fault driver process and passes to it the identifier of the faulting process and the page number. It unreadies the faulting process.
- A page-fault driver. This takes a message from the page-fault ISR and sends a message to the paging disk to retrieve the page whose reference caused the fault. If there are no free page frames in (main) physical store, it selects a victim page in physical store and sends it to the paging disk. When the faulting page is received from the paging disk, it is copied into a main store page whose physical address is identified with the logical page number in the faulting process. The faulting process is then readied and the driver waits for the next page fault.

The above scheme is sub-optimal. As part of the model, optimisations are suggested, particularly for the interactions between the driver and paging disk.

Even though it is sub-optimal, the above scheme is logically sufficient. It is, therefore, appropriate to concentrate on it as the model for this chapter. This exemplifies the method adopted in this book: capturing the logical functioning of the model is much more important than optimisation. The optimisation included here is introduced as an example of how it can be done without too much of a compromise to the model.

The structure of this kernel is shown in Figure 6.1. Comparison of this figure with the corresponding one in Chapter 4 (Figure 4.1) reveals their similarities. In the current kernel, virtual and not real storage forms the basis of the system. Apart from the need for structures and operations to support virtual storage (the subject of this chapter), the main difference lies in the kernel bootstrapping operations (which are not considered in this book).

## 6.3 Virtual Storage

In this section, the basic structures required to model a virtual store are introduced.

The following axiomatic definition defines the number of real pages (pages in real store or physical pages) and the size of the page frame. Neither constant is assigned a value, so the specification is of the loose variety.

$numrealpages : \mathbb{N}$
$framesize : \mathbb{N}$

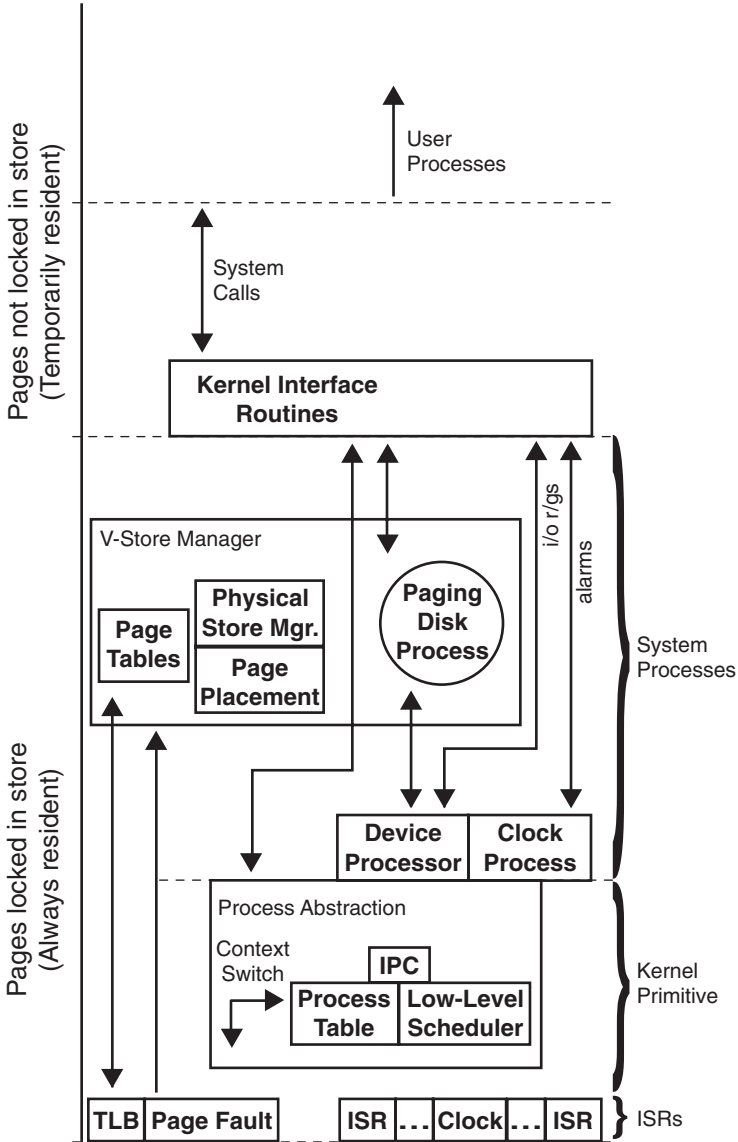The basic virtual address is represented by an atomic type:

Pages not locked in store
(Temporarily resident)

User
Processes

System
Calls

**Kernel Interface
Routines**

Pages locked in store
(Always resident)

V-Store Manager

**Page
Tables**

**Physical
Store Mgr.**

**Page
Placement**

**Paging
Disk
Process**

i/o r.'gs

alarms

System
Processes

**Device
Processor**

**Clock
Process**

Process Abstraction

Context
Switch

**IPC**

**Process
Table**

**Low-Level
Scheduler**

Kernel
Primitive

**TLB** **Page Fault**

**ISR** **...** **Clock** **...** **ISR**

ISRs

**Fig. 6.1.** *The layer-by-layer organisation of the kernel, including virtual storage-
management modules.*

[*VIRTUALADDRESS*]

> *maxvirtpagespersegment* : $\mathbb{N}$

This is the maximum number of virtual pages that a process can own.

$PAGEOFFSET == 1 .. framesize$
$PHYSICALPAGENO == 1 .. numrealpages$
$LOGICALPAGENO == \mathbb{N}$

where:

- *PAGEOFFSET* denotes the offsets into a page;
- *PHYSICALPAGENO* denotes the indices of pages in the main-store page frame;
- *LOGICALPAGENO* denotes the indices of logical pages (pages belonging to a process).

Note that *LOGICALPAGENO* is not bounded above. The reason for this is that the actual number of logical pages that a process can have is a hardware-dependent factor and is not relevant to the current exercise.

For every virtual address, the hardware performs an address translation that maps the virtual address to a logical frame number and an offset into that frame. The signature of this function, *addresstrans*, is:

> *addresstrans* : $VIRTUALADDRESS \rightarrow DECODEDADDRESS$

The definition of this function is hardware-specific and is, in any case, not particularly relevant to the current exercise.

The type *DECODEDADDRESS* is defined as:

$DECODEDADDRESS == LOGICALPAGENO \times FRAMEOFFSET$

and has the following projections:

> *dlogicalpage* : $DECODEDADDRESS \rightarrow LOGICALPAGENO$
> *dpageoffset* : $DECODEDADDRESS \rightarrow PAGEOFFSET$
> ─────────────
> $\forall \, addr : DECODEDADDRESS \bullet$
> $\quad dlogicalpage(addr) = fst \ \ addr$
> $\quad dpageoffset(addr) = snd \ \ addr$

For a segmented, paged architecture, the address-decoding function can be defined as:

> *saddresstrans* : $VIRTUALADDRESS \rightarrow SDECODEDADDRESS$

where:

$SDECODEDADDRESS == SEGMENT \times LOGICALPAGENO \times PAGEOFFSET$

and:

$saddrseg : SDECODEDADDRESS \rightarrow SEGMENT$
$spageno : SDECODEDADDRESS \rightarrow LOGICALPAGENO$
$spagoffset : SDECODEDADDRESS \rightarrow PAGEOFFSET$

$\forall saddr : SDECODEDADDRESS \bullet$
    $saddrseg(saddr) = fst\ saddr$
    $spageno(saddr) = fst(snd\ saddr)$
    $spagoffset(saddr) = snd^2\ saddr$

The *PAGEMAP* translates logical to physical page numbers. *PAGE-FRAME* translates physical page numbers to the actual pages in store. Finally, *PAGE* defines a storage page as a vector of *PSU*; the vector has *PAGEOFF-SET* elements. (*PSU* is, it will be recalled, the *Primary Storage Unit*, which is assumed to be a byte.)

$PAGEMAP == LOGICALPAGENO \nrightarrow PHYSICALPAGENO$
$PAGEFRAME == PHYSICALPAGENO \rightarrow PAGE$
$PAGE == PAGEOFFSET \rightarrow PSU$

The empty page is defined as follows:

$NullPage : PAGE$

$NullPage = (\lambda\, i : PAGEOFFSET \bullet 0)$

It is a page (vector of bytes), *PAGEOFFSET* bytes long, with each byte set to zero.

Pages are associated with a number of flags, including, among others:

- in core (i.e., in main store);
- locked In (i.e., must always remain in main store);
- shared (i.e., shared between at least two processes).

For each process, these properties can be represented by sets (thus simplifying the modelling of pages).

For the remainder of this section, a segmented virtual store will be modelled. The hardware nearest to the model presented here is the Intel X86 series. The reader is warned that this is a *logical* model of segmented, paged storage, not an exact rendition of any existing hardware implementations. Furthermore, details such as the *Translation Lookaside Buffer*, or *TLB*, (the associative store that typically holds a few entries from the current process' page table) are not modelled in detail, the reason being that, as usual, hardware differs considerably in implementation and, in particular, the size of the TLB can vary considerably among MMUs[1].

---

[1] *Memory Management Unit.*

Although most segmented hardware supports more than four segments per process, for the present, only three segments will be considered. Linux on X86 requires three segments: one each for code, stack and data, although the hardware permits a maximum of 16 segments (the virtual address size is 32 bits). The segment names are as follows:

$$SEGMENT == \{\mathsf{code}, \mathsf{data}, \mathsf{stack}, \mathsf{heap}, \ldots\}$$

A great many programming languages now require heap (dynamic) storage. A problem for dynamic store, as with stacks, is that, quite frequently, it has to be expanded. Within a three-segment organisation, the heap is part of either the stack or data segment; this can limit the maximum size of the heap somewhat on 32-bit machines. To simplify manipulation, it is assumed here that the data segment contains static data (global variables, literal pools, fixed-length buffers and so on) and the heap is given its own segment. The heap can, therefore, grow to the maximum segment size at runtime, as can the stack segment.

> $usedsegment : SEGMENT$
>
> ---
>
> $\forall\, s : SEGMENT \bullet$
> $\quad usedsegment(s) \Leftrightarrow s \in \{\mathsf{code}, \mathsf{data}, \mathsf{stack}, \mathsf{heap}\}$

It is now possible to define per-process page tables.

Each segment is composed of a number of pages. The following function translates a physical address and segment into a logical page number:

> $pages\_in\_segment : APREF \times SEGMENT \times$
> $\qquad (APREF \nrightarrow SEGMENT \nrightarrow \mathbb{F}\, LOGICALPAGENO) \rightarrow$
> $\qquad\qquad \mathbb{F}\, LOGICALPAGENO$
>
> ---
>
> $\forall\, p : APREF;\ sg : SEGMENT;\ f : APREF \nrightarrow SEGMENT \nrightarrow$
> $\qquad\qquad \mathbb{F}\, LOGICALPAGENO \bullet$
> $\quad pages\_in\_segment(p, sg, f) = f(p)(sg)$

The following (inverse) functions mark and unmark pages. They are higher-order functions that take the specification of a page and a page attribute map as arguments and return the modified page attribute map.

> $mark\_page : APREF \times SEGMENT \times LOGICALPAGENO \times$
> $\qquad (APREF \nrightarrow SEGMENT \nrightarrow \mathbb{F}\, LOGICALPAGENO) \rightarrow$
> $\qquad\qquad (APREF \nrightarrow SEGMENT \nrightarrow \mathbb{F}\, LOGICALPAGENO)$
> $unmark\_page : APREF \times SEGMENT \times LOGICALPAGENO \times$
> $\qquad (APREF \nrightarrow SEGMENT \nrightarrow \mathbb{F}\, LOGICALPAGENO) \rightarrow$
> $\qquad\qquad (APREF \nrightarrow SEGMENT \nrightarrow \mathbb{F}\, LOGICALPAGENO)$
>
> ---
>
> $\forall\, p : APREF;\ sg : SEGMENT;\ lpno : LOGICALPAGENO;$
> $\qquad f : APREF \nrightarrow SEGMENT \nrightarrow \mathbb{F}\, LOGICALPAGENO \bullet$
> $\quad mark\_page(p, sg, lpno, f) = f(p) \oplus \{sg \mapsto (f(p)(sg) \cup \{lpno\})\}$
> $\quad unmark\_page(p, sg, lpno, f) = f(p) \oplus \{sg \mapsto (f(p)(sg) \setminus \{lpno\})\}$

It is now proved that these two functions are mutual inverses.

**Proposition 124.** *mark_page and unmark_page$^{-1}$ are mutually inverse.*

PROOF. Write $f(p)(sg) = h$, then:

$$
\begin{aligned}
mark \;\;\;\; &= f(p) \oplus \{sg \mapsto (h \cup \{lpno\})\} \\
unmark &= f(p) \oplus \{sg \mapsto (h \setminus \{lpno\})\}
\end{aligned}
$$

Calculating the value of $unmark \circ mark$, we obtain:

$$ (f(p) \oplus \{sg \mapsto (h \cup \{lpno\})\}) \oplus \{sg \mapsto (h \setminus \{lpno\})\} $$

Writing $f(p)(sg) = s$, then $mark = s \cup \{lpno\}$ and $unmark = s \setminus \{lpno\}$, so:

$unmark \circ mark$
$$
\begin{aligned}
&= (s \setminus \{lpno\}) \circ (s \cup \{lpno\}) \\
&= (s \cup \{lpno\}) \setminus \{lpno\} \\
&= s
\end{aligned}
$$

Conversely:

$mark \circ unmark$
$$
\begin{aligned}
&= (s \cup \{lpno\}) \circ (s \setminus \{lpno\}) \\
&= (s \setminus \{lpno\}) \setminus \{lpno\} \\
&= s
\end{aligned}
$$

Therefore, *mark* and *unmark* are mutual inverses and the proposition is proved. □

The page table abstraction can be modelled as follows. The variable *free-pages* represents those pages in main store that are not allocated to any process. The page table proper is *pagetable*. The variables *executablepages*, *writablepages* and *readablepages* are intended to refer to pages the owner has marked executable (i.e., code pages), read-only (e.g., a constant data segment) and read-write (e.g., a stack). Pages can be shared between processes and some are locked into main store. When a page is locked, it cannot be removed from main store. The kernel's own storage is often marked as locked into main store. It is so locked because a page fault could prevent the kernel from responding in time to a circumstance. It is also necessary to keep track of those pages that are currently in main store: these are referred to as being "in core", hence the name of the variable, *incore*. The *pagecount* counts the number of pages in each segment of each process. There is an *a priori* limit to the number of pages in a segment and *pagecount* is intended to keep track of this and to provide a mechanism for raising an error condition if this limit is exceeded. The final variable, *smap*, is a relation between elements of *PAGE-SPEC*; it denotes those pages that are shared and it will be explained in more detail below.

There are different ways to organise page tables. The simplest is a linear sequence of page references. As virtual storage sizes increase, simple linear structures do not perform well, so tree-like structures are to be preferred. These trees can be arranged to perform mapping on two or three levels. The model defined here is intended to be suggestive of a tree structure, even though it can also be implemented as a table.

The class that follows defines an abstract data type. It represents the page table type. The type exports a large number of operations and has the most complex invariant in this book.

---

**PageTables**

$\uparrow$(*INIT*, *HaveFreePages*, *NumberOfFreePages*, *AllocateFreePage*,
    *MakePageFree*, *PhysicalPageNo*, *InitNewProcessPageTable*,
    *RemoveProcessFromPageTable*, *AddPageToProcess*, *HasPageInStore*,
    *IncProcessPageCount*, *DecProcessPageCount*,
    *LatestPageCount*, *UpdateMainstorePage*,
    *RemovePageFromProcessTable*, *RemovePageProperties*,
    *RemovePageFromProcess*, *IsPageInMainStore*, *MarkPageAsIn*,
    *MarkPageAsOut*, *IsSharedPage*, *MarkPageAsShared*,
    *UnsharePage*, *IsLockedPage*, *LockPage*, *UnlockPage*,
    *MakePageReadable*, *MakePageNotReadable*, *MakePageExecutable*,
    *IsPageExecutable*, *MakePageNotExecutable*, *MakePageWritable*,
    *IsPageWritable*, *MakePageNotWritable*)

---

*freepages* : $\mathbb{F}$ *PHYSICALPAGENO*
*pagetable* : *APREF* $\nrightarrow$ *SEGMENT* $\nrightarrow$ *PAGEMAP*
*executablepages*, *writablepages*, *readablepages*,
*sharedpages*, *lockedpages*,
*incore* : *APREF* $\nrightarrow$ *SEGMENT* $\nrightarrow$ $\mathbb{F}$ *LOGICALPAGENO*
*pagecount* : *APREF* $\nrightarrow$ *SEGMENT* $\nrightarrow$ $\mathbb{N}$
*smap* : *PAGESPEC* $\leftrightarrow$ *PAGESPEC*

---

*InvPageTables*

$0 \leq \#freepages \leq numrealpages$

dom *incore* $\subseteq$ dom *pagetable*

dom *sharedpages* $\subseteq$ dom *pagetable*

dom *lockedpages* $\subseteq$ dom *pagetable*

dom *pagecount* = dom *pagetable*

dom *executablepages* = dom *pagetable*

dom *writablepages* = dom *pagetable*

dom *readablepages* = dom *pagetable*

```
┌─ INIT ──────────────────────────────────────────────────────────
│  freepages' = ∅
│  dom pagetable' = ∅
│  dom sharedpages' = ∅
│  dom lockedpages' = ∅
│  dom incore' = ∅
│  dom pagecount' = ∅
│  dom executablepages' = ∅
│  dom writablepages' = ∅
│  dom readablepages' = ∅
│  dom smap = ∅
└─────────────────────────────────────────────────────────────────
```

$HaveFreePages \mathrel{\widehat{=}} \ldots$

$NumberOfFreePages \mathrel{\widehat{=}} \ldots$

$AllocateFreePage \mathrel{\widehat{=}} \ldots$

$MakePageFree \mathrel{\widehat{=}} \ldots$

$PhysicalPageNo \mathrel{\widehat{=}} \ldots$

$InitNewProcessPageTable \mathrel{\widehat{=}} \ldots$

$RemoveProcessFromPageTable \mathrel{\widehat{=}} \ldots$

$AddPageToProcess \mathrel{\widehat{=}} \ldots$

$HasPageInStore \mathrel{\widehat{=}} \ldots IncProcessPageCount \mathrel{\widehat{=}} \ldots$

$DecProcessPageCount \mathrel{\widehat{=}} \ldots$

$LatestPageCount \mathrel{\widehat{=}} \ldots$

$UpdateMainstorePage \mathrel{\widehat{=}} \ldots$

$RemovePageFromPageTable \mathrel{\widehat{=}} \ldots$

$RemovePageProperties \mathrel{\widehat{=}} \ldots$

$RemovePageFromProcess \mathrel{\widehat{=}} \ldots$

$IsPageInMainStore \mathrel{\widehat{=}} \ldots$

$MarkPageAsIn \mathrel{\widehat{=}} \ldots$

$MarkPageAsOut \mathrel{\widehat{=}} \ldots$

$IsSharedPage \mathrel{\widehat{=}} \ldots$

$MarkPageAsShared \mathrel{\widehat{=}} \ldots$

$UnsharePage \mathrel{\widehat{=}} \ldots$

$IsLockedPage \mathrel{\widehat{=}} \ldots$

$LockPage \mathrel{\widehat{=}} \ldots$

$UnlockPage \mathrel{\widehat{=}} \ldots$

$MakePageReadable \,\widehat{=}\, \ldots$

$MakePageNotReadable \,\widehat{=}\, \ldots$

$MakePageExecutable \,\widehat{=}\, \ldots$

$IsPageExecutable \,\widehat{=}\, \ldots$

$MakePageNotExecutable \,\widehat{=}\, \ldots$

$MakePageWritable \,\widehat{=}\, \ldots$

$IsPageWritable \,\widehat{=}\, \ldots$

$MakePageNotWritable \,\widehat{=}\, \ldots$

It will be noted that the invariant is partially stated in the class definition. The remainder is specified by the *InvPageTables* schema defined below after the other operations have been defined. This will bring the invariant closer to some of the proofs in which it is required.

The following schema represents the test that there are pages in main store (physical pages) that are free.

____ *HaveFreePages* _____
$freepages \neq \varnothing$

____ *NumberOfFreePages* _____
$np! : \mathbb{N}$
_____
$np! = \#freepages$

The following operation models the allocation of a free page to a process. It removes the page denoted by *ppno*! from the set of free pages, *freepages*.

____ *AllocateFreePage* _____
$\Delta(freepages)$
$ppno! : PHYSICALPAGENO$
_____
$ppno! \in freepages$
$freepages' = freepages \setminus \{ppno!\}$

**Proposition 125.** *AllocateFreePage implies that* $\#freepages' = freepages + 1$.

PROOF.

$\#freepages'$
$\quad = \#(freepages \setminus \{ppno?\})$
$\quad = \#freepages - \#\{ppno?\}$
$\quad = \#freepages - 1$

□

**Proposition 126.** *If freepages $= n$, AllocateFreePage$^n$ implies freepages$' =$ 0.*

PROOF.  By induction, using the last proposition.                    □

The next operation returns a page to the set of free pages.

---
*MakePageFree*
$\Delta(\textit{freepages})$
$ppno? : PHYSICALPAGENO$

---
$\textit{freepages}' = \textit{freepages} \cup \{ppno?\}$

---

**Proposition 127.** *MakePageFree implies that $\#\textit{freepages}' = \#\textit{freepages} - 1$.*

PROOF.

$\#\textit{freepages}'$
$\quad = \#(\textit{freepages} \cup \{ppno?\})$
$\quad = \#\textit{freepages} + \#\{ppno?\}$
$\quad = \#\textit{freepages} + 1$

                                                                          □

**Proposition 128.** *AllocateFreePage$[p/ppno!]$ $\,_9^\circ$ MakePageFree$[p/ppnp?]$ implies that freepages$' =$ freepages.*

PROOF.  The sequential composition can be written as:

$\exists \textit{freepages}'' : \mathbb{F}\, PHYSICALPAGENO \mid \textit{freepages}'' = \textit{freepages} \setminus \{ppno!\} \bullet$
$\quad \textit{freepages}' = \textit{freepages} \cup \{ppno?\}$

Renaming and simplifying:

$\textit{freepages}' = (\textit{freepages} \setminus \{p\}) \cup \{p\}$

Then:

$(\textit{freepages} \setminus \{p\}) \cup \{p\}$
$\quad = (\textit{freepages} \setminus \{p\}) \cup \{p\}$
$\quad = \textit{freepages} \cup (\{p\} \setminus \{p\})$
$\quad = \textit{freepages} \cup \varnothing$
$\quad = \textit{freepages}$

                                                                          □

The *PhysicalPageNo* operation contains a use of the *pagetable* variable. This variable is a higher-order function. Its use might appear a little odd.

Essentially, to obtain the physical page number corresponding to a logical page number, the process has to locate the segment in which the page occurs and then translate the logical page number.

---
*PhysicalPageNo*
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$
$ppgno! : PHYSICALPAGENO$

---

$ppgno! = pagetable(p?)(sg?)(lpno?)$

---

When a process is allocated, it is given an entry in the page table. The following schema models this operation. It just adds the process' identifier as a key in the subtables and sets everything to zero (empty or $\varnothing$).

---
*InitNewProcessPageTable*
$\Delta(incore, sharedpages, lockedpages)$
$p? : APREF$

---

$pagetable' = pagetable \cup \{p? \mapsto \mathsf{code} \mapsto \varnothing\} \cup \{p? \mapsto \mathsf{data} \mapsto \varnothing\}$
$\qquad \cup \{p? \mapsto \mathsf{stack} \mapsto \varnothing\} \cup \{p? \mapsto \mathsf{heap} \mapsto \varnothing\}$
$incore' = (incore \cup \{p? \mapsto \mathsf{code} \mapsto \varnothing\}) \cup (\{p? \mapsto \mathsf{data} \mapsto \varnothing\}$
$\qquad (\cup\{p? \mapsto \mathsf{stack} \mapsto \varnothing\}(\cup\{p? \mapsto \mathsf{heap} \mapsto \varnothing\})))$
$sharedpages' = (sharedpages \cup \{p? \mapsto \mathsf{code} \mapsto \varnothing\}) \cup (\{p? \mapsto \mathsf{data} \mapsto \varnothing\}$
$\qquad (\cup\{p? \mapsto \mathsf{stack} \mapsto \varnothing\}(\cup\{p? \mapsto \mathsf{heap} \mapsto \varnothing\})))$
$lockedpages' = (lockedpages \cup \{p? \mapsto \mathsf{code} \mapsto \varnothing\}) \cup (\{p? \mapsto \mathsf{data} \mapsto \varnothing\}$
$\qquad (\cup\{p? \mapsto \mathsf{stack} \mapsto \varnothing\}(\cup\{p? \mapsto \mathsf{heap} \mapsto \varnothing\})))$

---

**Proposition 129.** *InitNewProcessPageTable implies that the new process has no pages.*

PROOF.    For a process to have pages, it must have at least one page in *at least one* segment. However, for a process, $p$, and all segments, $sg$, $pagetable'(p)(sg) = \varnothing$.    □

**Corollary 11.** *InitNewProcessPageTable implies that the new process has no in-core pages.*

PROOF.  Similar to the above.    □

Similar results can be proved for all other page attributes, e.g., locked pages.

Conversely, when a process terminates or is killed, its storage is returned to the free pool and all of the information associated with it in the page tables is removed. The following schema models this operation:

```
┌─ RemoveProcessFromPageTable ──────────────────────────────
│  Δ(pagetable, incore, sharedpages, lockedpages)
│  p? : APREF
├───────────────────────────────────────────────────────────
│  pagetable' = {p?} ⊲ pagetable
│  incore' = {p?} ⊲ incore
│  sharedpages' = {p?} ⊲ sharedpages
│  lockedpages' = {p?} ⊲ lockedpages
└───────────────────────────────────────────────────────────
```

**Proposition 130.** *The predicate of RemoveProcessFromPageTable implies that $p \notin \text{dom}\, pagetable$.*

PROOF.  The first line of the predicate is $pagetable' = \{p?\} \triangleleft pagetable$. Taking domains:

$$
\begin{aligned}
\text{dom}\, pagetable' &= \text{dom}(\{p?\} \triangleleft pagetable) \\
&= (\text{dom}\, pagetable) \setminus \{p?\}
\end{aligned}
$$

$\square$

**Proposition 131.** *For any $p : APREF$, RemoveProcessFromPageTable$[p/p?]$ implies that the process:*

*1. is no longer incore, swappable, etc.; and*
*2. is no longer in the page table for its owning process.*

PROOF.  Each conjunct of the predicate employs the domain subtraction operation ($\triangleleft$) to remove $p$ from the domain of each function. This implies that $p$ is removed from each table.  $\square$

Propositions about page attributes can be proved. They follow the pattern of the last proposition.

When the storage image of a process is augmented by the addition of fresh pages, the following operation is the basic one used to extend the process' page table entry. Each page is specified as a process reference, a segment and a logical page number; in addition, the physical page number of the page to be added is also included. Since the process and segment are already present in the table, the logical to physical page number mapping is added to the table at the specified point.

```
┌─ AddPageToProcess ────────────────────────────────────────
│  Δ(pagetable)
│  p? : APREF
│  sg? : SEGMENT
│  lpno? : LOGICALPAGENO
│  ppno? : PHYSICALPAGENO
├───────────────────────────────────────────────────────────
│  pagetable' = pagetable(p?)(sg?) ∪ {lpno? ↦ ppno?}
└───────────────────────────────────────────────────────────
```

There now follows a predicate that returns true when the process specified by $p?$ has at least one page in main store:

```
┌─ HasPageInStore ──────────────────────────
│ p? : APREF
├───────────────────────────────────────────
│ p? ∈ dom incore
```

The per-segment page count is incremented by the following schema:

```
┌─ IncProcessPageCount ──────────────────────
│ Δ(pagecount)
│ p? : APREF
│ sg? : SEGMENT
├───────────────────────────────────────────
│ pagecount' = pagecount(p?) ⊕ {sg? ↦ pagecount(p?)(sg?) + 1}
```

The counter is decremented by the following schema:

```
┌─ DecProcessPageCount ──────────────────────
│ Δ(pagecount)
│ p? : APREF
├───────────────────────────────────────────
│ pagecount' = pagecount(p?) ⊕ {p? ↦ pagecount(p?)(sg?) − 1}
```

When a page is added to a segment, the page count is incremented. When a page is removed from a segment, the page count is decremented. The current value of the page count is obtained by the following schema:

```
┌─ LatestPageCount ──────────────────────────
│ p? : APREF
│ sg? : SEGMENT
│ lpno! : LOGICALPAGENO
├───────────────────────────────────────────
│ lpno! = pagecount(p?)(sg?)
```

If the logical to physical page mapping is changed, the following schema performs the update in the page table.

```
┌─ UpdateMainStorePage ──────────────────────
│ Δ(pagetable)
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
│ ppno? : PHYSICALPAGENO
├───────────────────────────────────────────
│ pagetable' = pagetable(p?)(sg?) ⊕ {lpno? ↦ ppno?}
```

When a page is removed from the page table, the entry representing it must be removed. The removal operation is defined as follows:

$\underline{\quad RemovePageFromPageTable \quad}$
$\Delta(pagetable)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$
$\overline{\qquad}$
$(\exists\, pmap : PAGEMAP \mid pmap = pagetable(p?)(sg?) \bullet$
$\qquad pagetable' = pagetable(p?) \oplus \{sg? \mapsto (\{lpno?\} \lhd pmap)\})$

The removal of a page also requires the removal of the attributes of that page. The attributes are removed using the *unmark_page* function (when a page is allocated, the attributes it possesses are marked using the *mark_page* function).

$\underline{\quad RemovePageProperties \quad}$
$\Delta(executablepages, readablepages, writablepages, sharedpages,$
$\qquad lockedpages, incore)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$
$\overline{\qquad}$
$executablepages' = unmark\_page(p?, sg?, lpno?, executablepages)$
$readablepages' = unmark\_page(p?, sg?, lpno?, readablepages)$
$writablepages' = unmark\_page(p?, sg?, lpno?, writablepages)$
$sharedpages' = unmark\_page(p?, sg?, lpno?, sharedpages)$
$lockedpages' = unmark\_page(p?, sg?, lpno?, lockedpages)$
$incore' = unmark\_page(p?, sg?, lpno?, incore)$

Finally, the high-level operation to remove a page from a process is defined as follows:

$RemovePageFromProcess \,\widehat{=}$
$\qquad RemovePageFromProcessTable \,\mathbin{{}_9^\circ}$
$\qquad MakePageFree \,\mathbin{{}_9^\circ}$
$\qquad RemovePageProperties$

It is possible to determine whether a page is in main store by determining whether it is in the *incore* attribute. The following schema defines this predicate. Note that it uses the *pages_in_segment* function.

$\underline{\quad IsPageInMainStore \quad}$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$
$\overline{\qquad}$
$lpno? \in pages\_in\_segment(p?, sg?, incore)$

**Proposition 132.** *IsPageInMainStore iff lpno? is an in-core page.*

PROOF. The predicate states that $lpno? \in pages\_in\_segment(p?, sg?, incore)$. By the definition of $pages\_in\_segment$,

$$lpno? \in incore(p?)(sg?)$$

□

**Proposition 133.** *IsSharedPage iff lpno? is a shared page; that is, iff lpno? is an element of sharedpages(p)(sg), for some p and sg.*

PROOF. Similar to the previous proof.    □

**Proposition 134.** *IsLockedPage iff lpno? is a locked page; that is, iff lpno? is an element of lockedpages.*

PROOF. Similar to the previous proof.    □

When a page is swapped into main store, it is marked as being "in". The following schema performs this marking. The schema that immediately follows marks pages as "out" (i.e., as not being main-store resident).

```
__ MarkPageAsIn _____
  Δ(incore)
  p? : APREF
  sg? : SEGMENT
  lpno? : LOGICALPAGENO
 ────────────────────────
  incore' = mark_page(p?, sg?, lpno?, incore)
```

**Proposition 135.** *MarkPageAsIn implies that lpno? is an element of incore'.*

PROOF. As can be seen from the schema, the predicate is $incore' = mark\_page(p?, sg?, lpno?, incore)$. Substituting the definition of $mark\_page$:

$incore' = mark\_page(p?, sg?, lpno?, incore)$
$\quad = incore(p?) \oplus \{sg \mapsto (incore(p?)(sg) \cup \{lpno?\})\}$

□

**Proposition 136.** *For fixed arguments, $p : APREF$, $s : SEGMENT$ and $l : LOGICALPAGENO$, $MarkPageAsIn[p/p?, s/sg?, l/lpno?]^n$ has the same effect as $MarkPageAsIn[p/p?, s/sg?, l/lpno?]$.*

PROOF.  This proposition is to be taken as:

$$MarkPageAsIn^n \Leftrightarrow MarkPageAsIn$$

The proposition is proved by substitution from the following general property of sets:

$$(S \cup \{x\}) \cup \{x\} = S \cup \{x\}$$

(i.e., the absorbtive law of set union).                                         □

---
*MarkPageAsOut* ────────────────────────────────

$\Delta(incore)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

─────

$incore' = unmark\_page(p?, sg?, lpno?, incore)$

---

**Proposition 137.** *MarkPageAsOut is satisfied iff lpno? is not an element of incore.*

PROOF.  Substituting the definition of *unmark_page* into the predicate of the schema *MarkPageAsOut*:

$incore' = unmark\_page(p?, sg?, lpno?, incore)$
$\quad = incore(p?) \oplus \{sg? \mapsto (incore(p?)(sg?) \setminus \{lpno?\})\}$

□

**Proposition 138.** *For fixed arguments, $p : APREF$, $s : SEGMENT$ and $l : LOGICALPAGENO$, $MarkPageAsOut[p/p?, s/sg?, l/lpno?]^n$ has the same effect as $MarkPageAsOut[p/p?, s/sg?, l/lpno?]$.*

PROOF.  The statement of the proposition is to be taken as:

$$MarkPageAsOut^n \Leftrightarrow MarkPageAsOut$$

The proposition is proved by substitution from the following general property of sets:

$$(S \setminus \{x\}) \setminus \{x\} = S \setminus \{x\}$$

□

**Proposition 139.** $MarkPageAsIn[l/lpno?] \,^\circ_9 MarkPageAsOut[l/lpno?]$ *implies that $incore' = incore$.*

PROOF. Writing out the predicates and performing the obvious substitutions:

$$unmark\_page(p?, sg?, lpno?, mark\_page(p?, sg?, lpno?, incore))$$

The result follows from the fact that *unmark_page* and *mark_page* are mutually inverse.                                                                    □

**Proposition 140.** *MarkPageAsOut[l/lpno?]⨾MarkPageAsIn[l/lpno?] implies that incore = incore′.*

PROOF. Writing out the predicates and performing the obvious substitutions:

$$mark\_page(p?, sg?, lpno?, unmark\_page(p?, sg?, lpno?, incore))$$

The result follows from the fact that *unmark_page* and *mark_page* are mutually inverse.                                                                    □

The next few schemata set and unset attributes in pages. The attributes are represented by the various tables in the *PageTables* class, such as *sharedpages*, *readable* and *locked*. The schemata naturally fall into three sets: one to perform a test and one to set the attribute and one to unset it. The schemata in each of these sets have the same structure. That structure is the obvious one and is quite simple. For these reasons, the schemata will not be described in English: the formal notation can stand on its own.

---

__ *IsSharedPage* _____

p? : APREF
sg? : SEGMENT
lpno? : LOGICALPAGENO
_____

lpno? ∈ pages_in_segment(p?, sg?, sharedpages)

---

__ *MarkPageAsShared* _____

Δ(sharedpages)
p? : APREF
sg? : SEGMENT
lpno? : LOGICALPAGENO
_____

sharedpages′ = mark_page(p?, sg?, lpno?, sharedpages)

---

__ *UnsharePage* _____

Δ(sharedpages)
p? : APREF
sg? : SEGMENT
lpno? : LOGICALPAGENO
_____

sharedpages′ = unmark_page(p?, sg?, lpno?, sharedpages)

---

___ *IsLockedPage* _____

$p? : APREF$

$sg? : SEGMENT$

$lpno? : LOGICALPAGENO$

_____

$lpno? \in pages\_in\_segment(p?, sg?, lockedpages)$

___ *LockPage* _____

$\Delta(lockedpages)$

$p? : APREF$

$sg? : SEGMENT$

$lpno? : LOGICALPAGENO$

_____

$lockedpages' = mark\_page(p?, sg?, lpno?, lockedpages)$

___ *UnlockPage* _____

$\Delta(lockedpages)$

$p? : APREF$

$sg? : SEGMENT$

$lpno? : LOGICALPAGENO$

_____

$lockedpages' = unmark\_page(p?, sg?, lpno?, lockedpages)$

___ *MakePageReadable* _____

$\Delta(readablepages)$

$p? : APREF$

$sg? : SEGMENT$

$lpno? : LOGICALPAGENO$

_____

$readablepages' = mark\_page(p?, sg?, lpno?, readablepages)$

___ *IsPageReadable* _____

$p? : APREF$

$sg? : SEGMENT$

$lpno? : LOGICALPAGENO$

_____

$lpno? \in readablepages(p?)(sg?)$

___ *MakePageNotReadable* _____

$\Delta(readablepages)$

$p? : APREF$

$sg? : SEGMENT$

$lpno? : LOGICALPAGENO$

_____

$readablepages' = unmark\_page(p?, sg?, lpno?, readablepages)$

___ *MakePageExecutable* _____

$\Delta(executablepages)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

_____

$executablepages' = mark\_page(p?, sg?, lpno, executablepages)$

___ *IsPageExecutable* _____

$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

_____

$lpno? \in executablepages(p?)(sg?)$

___ *MakePageNotExecutable* _____

$\Delta(executablepages)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

_____

$executablepages' = unmark\_page(p?, sg?, lpno, executablepages)$

___ *MakePageWritable* _____

$\Delta(writablepages)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

_____

$writablepages' = mark\_page(p?, sg?, lpno, writablepages)$

___ *IsPageWritable* _____

$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

_____

$lpno? \in writablepages(p?)(sg?)$

___ *MakePageNotWritable* _____

$\Delta(writablepages)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

_____

$writablepages' = unmark\_page(p?, sg?, lpno, writablepages)$

Finally, we come to *InvPageTables*, the page table invariant. As can be seen, this is quite an invariant. Because of its length, it was not possible to make it fit within the Object-Z class box without obfuscating the specification even more. The invariant is mostly concerned with ensuring that each page is represented correctly. For example, no free page has a corresponding page-table entry; all locked pages are always in main store, and so on. Readers can inspect the various clauses of the invariant for themselves.

This invariant clearly demonstrates the need for higher-order logics!

$InvPageTables \cong$
$(\forall\, ppno : PHYSICALPAGENO \bullet$
$\quad ppno \in freepages \Leftrightarrow$
$\quad\quad \neg\,(\exists\, p : APREF;\; sg : SEGMENT;\; lpno : LOGICALPAGENO \bullet$
$\quad\quad\quad\quad ppno = pagetables(p)(sg)(lpno)))$
$(\forall\, p : APREF;\; sg : SEGMENT \bullet$
$\quad p \in \mathrm{dom}\, lockedpages \wedge sg \in \mathrm{dom}\, lockedpages(p) \Rightarrow$
$\quad\quad lockedpages(p)(sg) \subseteq incore(p)(sg))$
$(\forall\, p : APREF;\; sg : SEGMENT \bullet$
$\quad p \in \mathrm{dom}\, incore \wedge sg \in \mathrm{dom}\, incore(p) \Rightarrow$
$\quad\quad incore(p)(sg) \subseteq \mathrm{dom}\, pagetable(p)(sg))$
$(\forall\, p : APREF;\; sg : SEGMENT \bullet$
$\quad p \in \mathrm{dom}\, sharedpages \wedge sg \in \mathrm{dom}\, sharedpages(p) \Rightarrow$
$\quad\quad sharedpages(p)(sg) \subseteq \mathrm{dom}\, pagetable(p)(sg))$
$(\forall\, p : APREF;\; sg : SEGMENT \bullet$
$\quad p \in \mathrm{dom}\, lockedpages \wedge sg \in \mathrm{dom}\, lockedpages(p) \Rightarrow$
$\quad\quad lockedpages(p)(sg) \subseteq \mathrm{dom}\, pagetable(p)(sg))$

These clauses are intended to represent the disjointness of segments.

$(\forall\, p : APREF \mid p \in \mathrm{dom}\, pagetable \bullet$
$\quad (\forall\, sg_1, sg_2 : SEGMENT \bullet$
$\quad\quad (sg_1 \neq sg_2 \wedge$
$\quad\quad\quad sg_1 \in \mathrm{dom}\, pagetable(p) \wedge sg_2 \in \mathrm{dom}\, pagetable(p)) \Rightarrow$
$\quad\quad\quad\quad \mathrm{dom}\, pagetable(p)(sg_1) \cap \mathrm{dom}\, pagetable(p)(sg_2) = \varnothing))$
$(\forall\, p : APREF \mid p \in \mathrm{dom}\, pagetable \bullet$
$\quad (\forall\, sg_1, sg_2 : SEGMENT \bullet$
$\quad\quad (sg_1 \neq sg_2 \wedge$
$\quad\quad\quad sg_1 \in \mathrm{dom}\, pagetable(p) \wedge sg_2 \in \mathrm{dom}\, pagetable(p)) \Rightarrow$
$\quad\quad\quad\quad \mathrm{ran}\, pagetable(p)(sg_1) \cap \mathrm{ran}\, pagetable(p)(sg_2) = \varnothing))$
$(\forall\, p : APREF;\; sg : SEGMENT \bullet$
$\quad p \in \mathrm{dom}\, executablepages \wedge sg \in executablepages(p) \Rightarrow$
$\quad\quad executablepages(p)(sg) \subseteq \mathrm{dom}\, pagetable(p)(sg))$
$(\forall\, p : APREF;\; sg : SEGMENT \bullet$
$\quad p \in \mathrm{dom}\, readablepages \wedge sg \in readtablepages(p) \Rightarrow$
$\quad\quad readablepages(p)(sg) \subseteq \mathrm{dom}\, pagetable(p)(sg))$

$$(\forall\, p : APREF;\ sg : SEGMENT \bullet$$
$$\quad p \in \mathrm{dom}\ writablepages \wedge sg \in writablepages(p) \Rightarrow$$
$$\qquad writablepages(p)(sg) \subseteq \mathrm{dom}\ pagetable(p)(sg))$$


$$(\forall\, pg : PAGESPEC \mid pg \in \mathrm{dom}\ pmap \bullet$$
$$\quad (\exists\, p : APREF;\ sg : SEGMENT;\ lpno : LOGICALPAGENO \bullet$$
$$\qquad (p = pgspecpref(pg)\ \wedge$$
$$\qquad\quad sg = pgspecseg(pg)\ \wedge$$
$$\qquad\quad lpno = pgspeclpno(pg)) \Rightarrow$$
$$\qquad (p \in \mathrm{dom}\ pagetable \wedge sg \in \mathrm{dom}\ pagetable(p)\ \wedge$$
$$\qquad\quad lpno \in \mathrm{dom}\ pagetable(p)(sg))))$$
$$(\forall\, pg : PAGESPEC \mid pg \in \mathrm{ran}\ pmap \bullet$$
$$\quad (\exists\, p : APREF;\ sg : SEGMENT;\ lpno : LOGICALPAGENO \bullet$$
$$\qquad (p = pgspecpref(pg)\ \wedge$$
$$\qquad\quad sg = pgspecseg(pg)\ \wedge$$
$$\qquad\quad lpno = pgspeclpno(pg)) \Rightarrow$$
$$\qquad (p \in \mathrm{dom}\ pagetable \wedge sg \in \mathrm{dom}\ pagetable(p)\ \wedge$$
$$\qquad\quad lpno \in \mathrm{dom}\ pagetable(p)(sg))))$$

**Proposition 141.**

$$InitNewProcessPageTable \Rightarrow$$
$$\quad (\forall\, s : SEGMENT \bullet$$
$$\qquad incore = \varnothing$$
$$\qquad \wedge\ sharedpages = \varnothing$$
$$\qquad \wedge\ lockedpages = \varnothing)$$

PROOF.  By the predicate of the *InitNewProcessPageTable* predicate.  □


**Proposition 142.** $\forall\, p : PAGE \bullet locked(p) \Leftrightarrow \neg\ swappable(p)$

PROOF.  The class invariant states that dom *lockedpages* $\subseteq$ dom *pagetables*. This permits us to infer that:

$$\forall\, p : LOGICALPAGENO \bullet$$
$$\quad p \in lockedpages(p)(sg)$$
$$\qquad \Rightarrow p \in (\mathrm{dom}\ pagetable(p)(sg))$$

for all $p \in IPREF$ and $sg \in SEGMENT$.
    Again, by the same invariant:

$$lockedpages(p)(sg) \subseteq incore(p)(sg)$$

(again, for all $p$ and $sg$ as above).
    These two formulæ ensure that every locked page exists in store.
$\Rightarrow$: By the predicate of *PageFaultDriver.findVictimLogicalPage* (q.v.):

$pagetable(p)(sg)(l) = v! \land l \notin lockedpage$

First, applying the substitutions $[p/p!, sg/sg!, l/l!]$, then simplifying and rearranging the predicate of *findVictimLogicalPage*, we obtain:

$p \in \text{dom } pagetable$
$\quad \land \, sg \in \text{dom } pagetable(p)$
$\quad \land \, l \in \text{dom } pagetable(p)(sg)$
$\quad \land \, l \notin lockedpages(p)(sg)$
$\quad \land \, pagetable(p)(sg)(l) = v$

where $v$ is the page to be swapped.

From this, it can be concluded that $v$ cannot be a locked page.

$\Leftarrow$. If a page is not swappable, it is locked. Consider the set of swappable pages, $S$, by the definition of *findVictimLogicalPage* in the substitution instance above, any logical page, $l$, cannot be in $S$ because $l \notin lockedpages(p)(sg)$. $\quad\square$

**Proposition 143.** *For all processes, lockedpages $\cap$ swappablepages $= \varnothing$.*

PROOF. This follows immediately from the previous proposition. $\quad\square$

**Proposition 144.** *If a page, p, is in main store, IsPageInMainStore[p/...] is satisfied.*

PROOF. By the predicate of *IsPageInMainStore*, it is clear that

$$lpno? \in pages\_in\_segment(p, s, incore)$$

and that: $pages\_in\_segment(x, y, f) = f(x)(y)$; its range is a set.

For a page actually to be in store, the following must be satisfied:

*IsPageInMainstore* implies that:
$\quad (\exists\, p : IPREF;\ sg : SEGMENT;\ pp : PHYSICALPAGENO \bullet$
$\quad\quad pagetable(p)(sg)(lpno) = pp$
$\quad\quad \land\, pp \notin freepages)$

By the invariant:

$\forall\, pp : PHYSICALPAGENO \bullet$
$\quad pp \in freepages \Rightarrow$
$\quad\quad \neg\,(\exists\, p : IPREF;\ sg : SEGMENT;\ l : LOGICALPAGENO \bullet$
$\quad\quad\quad pp = pagetables(p)(sg)(l)$

By the application of contraposition $(p \Rightarrow q \Leftrightarrow q \Rightarrow p)$, the result is obtained.

It is now necessary to show that

$l \in incore(p)(sg) \Rightarrow$
$\quad (\exists\, pp_1 : PHYSICALPAGENO \bullet pagetables(p)(sg)(l) = pp)$

This can be done using the invariant. $\quad\square$

**Proposition 145.** *If a page, p, is locked, IsPageInMainStore[p/...] is satisfied.*

PROOF.  Similar to that of the last proof.                                    □

**Proposition 146.** *A locked page can never be swapped out.*

PROOF.  By Propositions 142 and 144.                                          □

**Proposition 147.** *A free page is in the process table of no process.*

PROOF.  Assume that process, $p$, has segment, $s$, in which the logical page number, $l$ is mapped to physical page, $n$. In this case, $n = pagetable(p)(s)(l)$.
By the class invariant,

$\forall\, ppno : PHYSICALPAGENO \bullet$
  $ppno \in freepages \Leftrightarrow$
      $\neg\, (\exists\, pid : APREF;\ sg : SEGMENT;\ lpn : LOGICALPAGENO \bullet$
        $ppno = pagetable(pid)(sg)(lpn))$

By Universal Instantiation:

$n \in freepages \Leftrightarrow$
  $\neg\, (\exists\, pid : APREF;\ sg : SEGMENT;\ lpn : LOGICALPAGENO \bullet$
        $n = pagetable(pid)(sg)(lpn))$

From the assumption that $n \in freepages$, a contradiction ensues. Therefore,

$$n \notin freepages$$

as required.                                                                 □

**Proposition 148.** *A free page is not incore.*

PROOF.  The invariant states that:

$\forall\, ppno : PHYSICALPAGENO \bullet$
  $ppno \in freepages \Leftrightarrow$
      $\neg\, (\exists\, pid : APREF;\ sg : SEGMENT;\ lpn : LOGICALPAGENO \bullet$
        $ppno = pagetable(pid)(sg)(lpn))$

This implies that:

$$freepages \cap pagetable(p)(sg) = \varnothing \qquad (6.1)$$

The class invariant also states that:

$p \in \mathrm{dom}\, incore \land sg \in \mathrm{dom}\, incore(p) \Rightarrow$
$$incore(p)(sg) \subseteq \mathrm{dom}\, pagetable(p)(sg)$$

This and equation (6.1) above imply that $pg \in freepages \Leftrightarrow \neg\, incore(p)(sg)(l)$ for all values of $p$, $sg$ and $l$.                                                                  □

**Proposition 149.** *A free page is not:*

1. *executable;*
2. *readable; or*
3. *writable.*

PROOF.  Since a free page is not in the page table, it follows from the invariant that it cannot have any of these attributes.                                                             □

### 6.3.1 The Paging Disk Process

The paging disk holds pages while they are not in main store. The virtual store software copies pages to and from the paging disk to implement the swapping process that underlies the huge address space illusion. Pages are swapped out of store when they are not required and swapped back in when their owning process (or one of them, if the page is shared) refers to them.

The paging disk is part of the subsystem whose design will be discussed in the next few subsections. The subsystem's organisation and interactions are shown in Figure 6.2. It consists of an ISR, a handler or driver process and the paging disk proper.

The paging disk is assumed to be infinite in size. It is represented by a mapping from logical to physical pages (*pagemap*). In addition, the disk has an interface to the message-passing system (*msgmgr*) so that it can communicate with the other processes in the storage-management subsystem. For most of the time, *pagemap* will be the focus of attention.

The class defining the paging disk process is now defined.

---
*PagingDiskProcess*
$\lceil(INIT, PageIsOnDisk, StorePageOnDisk, RetrievePageFromDisk,$
     $RemoveProcessFromPagingDisk, OnPageRequest)$

---
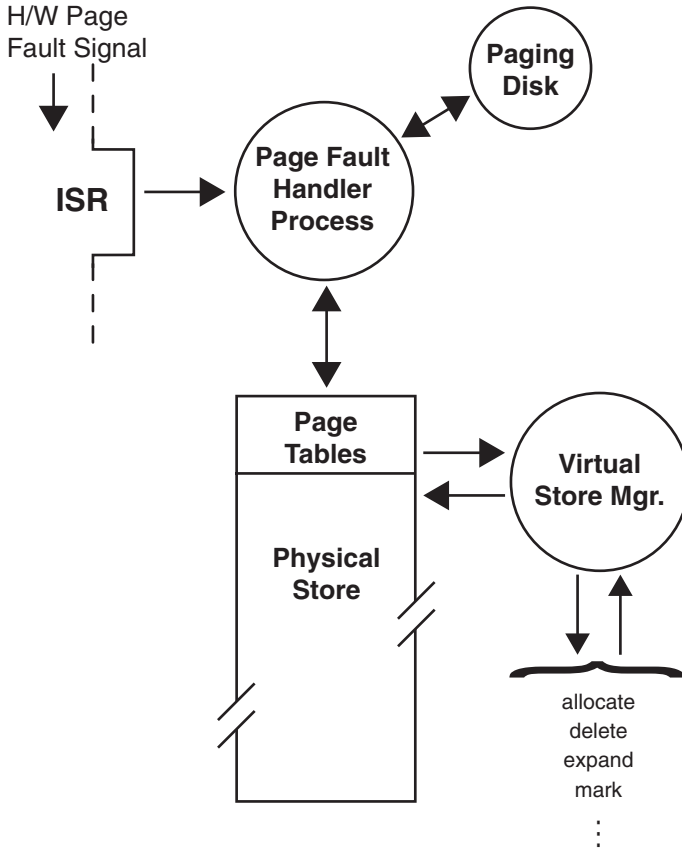$pagemap : APREF \nrightarrow SEGMENT \nrightarrow LOGICALPAGENO \nrightarrow PAGE$
$msgmgr : MsgMgr$

**Fig. 6.2.** *Interactions between virtual storage components.*

---

$\underline{\quad INIT \underline{\phantom{aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa}}}$
$msgs? : MsgMgr$
___
$msgmgr' = msgs?$
$\mathrm{dom}\, pagemap' = \varnothing$

$PageIsOnDisk \;\hat{=}\; \ldots$

$StorePageOnDisk \;\hat{=}\; \ldots$

$RetrievePageFromDisk \;\hat{=}\; \ldots$

$RemoveProcessFromPagingDisk \;\hat{=}\; \ldots$

$OnPageRequest \;\hat{=}\; \ldots$

The following schema is a predicate that is true when the specified page of the specified process is present on the paging disk.

```
┌─ PageIsOnDisk ──────────────────────────────────────────
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
├──────────────────────────────────────────────────────────
│ p? ∈ dom pagemap
│ sg? ∈ dom pagemap(p?)
│ lpno? ∈ dom pagemap(p?)(sg?)
└──────────────────────────────────────────────────────────
```

Pages are stored on the paging disk when they are swapped out of main store. In order for a page to be stored on the paging disk, it must be placed there and indexed properly. The *StorePageOnDisk* operation does this:

```
┌─ StorePageOnDisk ──────────────────────────────────────
│ Δ(pagemap)
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
│ pg? : PAGE
├──────────────────────────────────────────────────────────
│ pagemap' = pagemap ⊕ {p? ↦ {sg? ↦ {lpno? ↦ pg?}}}
└──────────────────────────────────────────────────────────
```

**Proposition 150.** *If a logical page image is already on disk and that page is swapped out again, that page image is overwritten.*

PROOF. Let $p = pagemap(p?)(sg?)(lpno?)$ and let $pg?$ be $pagemap'(p?)(sg?)$ $(lpno?)$, so $\{p? \mapsto \{sg? \mapsto \{lpno? \mapsto p\}\}\} \in pagemap$ and $\{p? \mapsto \{sg? \mapsto \{lpno? \mapsto pg?\}\}\} \in pagemap'$. Therefore, $pagemap' = pagemap \oplus \{p? \mapsto \{sg? \mapsto \{lpno? \mapsto pg?\}\}\}$.                                  □

When its owning process references a page that is not in main store, the paging disk is instructed to retrieve it from the paging disk. The following schema models the basics of this operation:

```
┌─ RetrievePageFromDisk ─────────────────────────────────
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
│ pg! : PAGE
├──────────────────────────────────────────────────────────
│ pg! = pagemap(p?)(sg?)(lpno?)
└──────────────────────────────────────────────────────────
```

When a process terminates (or is terminated), all of its pages must be removed from the paging disk. In our model, this can be done very easily by

removing the process' identifier from the domain of the *pagemap*. This removes all references to the process from the map—hence, removes the process from the disk, thus:

```
┌─ RemoveProcessFromPagingDisk ──────────────────────────────
│ Δ(pagemap)
│ p? : APREF
│ ──────────────────────
│ pagemap' = {p?} ◁ pagemap
└────────────────────────────────────────────────────────────
```

**Proposition 151.** *The predicate of RemoveProcessFromPagingDisk implies that $p \notin \mathrm{dom}\,pagemap'$.*

PROOF.

$\mathrm{dom}\,pagemap'$
$\quad = \mathrm{dom}(\{p?\} ◁ pagemap)$
$\quad = (\mathrm{dom}\,pagemap) \setminus \{p?\}$

Therefore, $p \notin \mathrm{dom}\,pagemap'$.                                    □

The paging disk has an operation that handles requests for operations. The operations it can perform are:

- store a page (*STOPG*);
- retrieve a page (*GTPG*), and
- delete a process from the paging disk (*DELPROCPG*).

The following schema defines this operation. The schema uses overloaded calls to the *RcvMessage* primitive provided by the message-handling subsystem. The operation is really just a large "or" based on the type of message that has just been received. The infinite loop is modelled by the exterior universal quantifier. While the operation has nothing to do, it waits for the next message to arrive.

```
┌─ OnPageRequest ────────────────────────────────────────────
│ ∀ i : 1 .. ∞ •
│     (∃ pdsk, fhandler : APREF; msg : MSG |
│                   pdsk = pagedisk ∧ fhandler = faultdrvr •
│           msgmgr.RcvMessage[pdsk/caller?, fhandler/src?, msg/m!]
│           ∧ (∃ p : APREF; sg : SEGMENT;
│                       lpno : LOGICALPAGENO; pg : PAGE •
│               (msg = STOPG⟨⟨p, sg, lpno, pg⟩⟩
│                     ∧ storePageOnDisk[p/p?, sg/sg?, lpno/lpno?, pg/pg?])
└────────────────────────────────────────────────────────────
```

$$\lor \ (msg = GTPG\langle\!\langle p, sg, lpno \rangle\!\rangle$$
$$\land \ retrievePageFromDisk[p/p?, sg/sg?, lpno/lpno?, pg/pg!]$$
$$\land \ (\exists \ rpmsg : MSG \mid rpmsg = ISPG\langle\!\langle p, sg, lpno, pg \rangle\!\rangle \bullet$$
$$msgmgr.SendMessage$$
$$[fhandler/dest?, pdsk/src?, rpmsg/m?]))$$
$$\lor \ (msg = DELPROCPG\langle\!\langle p \rangle\!\rangle$$
$$\land \ removeProcessFromPagingDisk[p/p?])))$$

**Proposition 152.** *OnPageRequest does what it should.*

PROOF. The proof is by cases on message type. It is assumed that all domains are correct so the error cases ignored in the schema need not be introduced here.

Case 1. $m = STOPG$, a request to store a page on disk. The predicate of $storePageOnDisk[p/p?, sg/sg?, lpno/lpno?, pg/pg?]$ states:

$$pagemap' = pagemap \oplus \{p \mapsto \{sg \mapsto \{lpno \mapsto pg\}\}\}$$

so $pg? = pagemap'(p)(sg)(lpno)$.

Case 2. $m = GTPG$, a request to retrieve a page. This follows from the fact that $retrievePageOnDisk$ is a function: $pg? = pagemap(p?)(sg?)(lpno?)$.

Case 3. $m = DELPROCPG$. By the last proposition. $\qquad\square$

### 6.3.2 Placement: Demand Paging and LRU

There are significant issues to be resolved in the design of the virtual storage mechanism. This section deals with the general area of *placement*. Placement is concerned with where pages are to be removed and included in main store. When a process refers to a page that is not currently in main store, it generates a *page fault*, an asynchronous signal that interrupts the process and leads to the satisfaction of the reference. To do this, the support software has to identify the page that is being referenced and then identify a page in physical store that can be swapped out to the paging disk, thus making space for the referenced page to be copied into main store. The issue is slightly complicated by the fact that the system might have some free physical pages in main store (indeed, it might be a policy to keep a block of such pages in reserve). For present purposes, it is assumed that *all* physical pages are allocated and that there is no pool of free pages kept in reserve.

The placement algorithm just outlined is *demand paging*. This is the most common approach. It is documented, like many other possible approaches, in most textbooks on operating systems (e.g., [26, 11, 29, 5]). Demand paging is the most commonly used approach and makes reasonable assumptions about the hardware and the software. It just assumes that the hardware can detect

a page fault and that the software can find the referenced page and locate the page in main store where it can be stored; if there is no free page, demand paging assumes that there is a way to find a victim page that can be swapped out to make space.
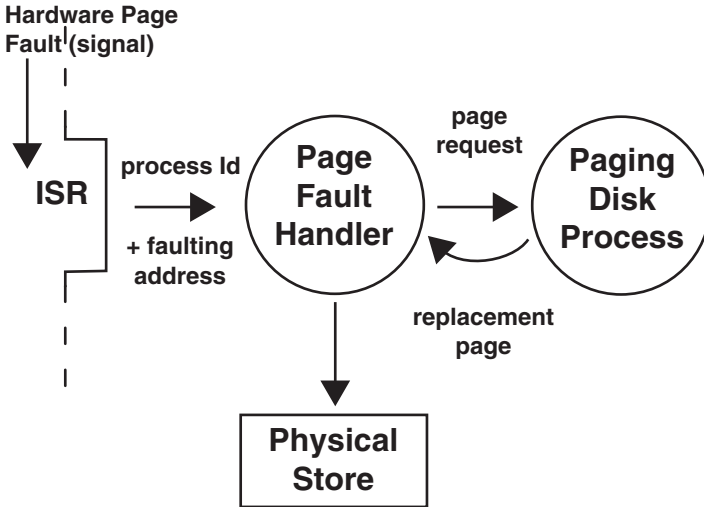


**Fig. 6.3.** *Process organisation for handling page faults.*

The placement algorithm used in this model has two identifiable aspects:

1. finding pages to remove;
2. pages to include.

The second aspect is solved for us: the pages to include are always those that cause a page fault when referenced by a process.

The general organisation of the processes that handle page faults is shown in Figure 6.3.

### 6.3.3 On Page Fault

A page fault occurs when a virtual address refers to a page that is not in physical store.

When such a fault occurs, it is assumed that the address causing the fault is available to the operating system. From this address, it is assumed that the following parameters can be computed: segment number, logical page number and the offset into the logical page. In addition, it is assumed that the

identity of the process causing the page fault can be determined (it should be the currently running process on a uni-processor; in the terms of this book, the value of *currentp*).

The case in which a page is shared is identical. The important thing to remember is that the owning process causes the page fault. If a page is shared, the important thing is for the page not to be swapped in (or out) more than once. That is why the *smap* is included in *PageTables*. When a shared page is to be swapped, this map is consulted. If the page is shared, it should already be in store. Conversely, if a page is to be swapped in, the *smap* should be consulted for the same reason.

In this section, a page that is to be removed from main store in order to free a page frame will be called a *victim* or *candidate*. It is always a restriction on victim page selection that victim pages are never locked into main store.

To define a relatively simple candidate-finding algorithm, it is necessary to associate page frames in main store with a bit that is set by the hardware whenever a page is referenced and a counter. The counter is a single byte whose value is computed by rotating the reference bit into the top bit of the counter byte and or-ing the counter with the contents of the counter shifted down one bit (ignoring the bottom bit). The types of the reference bit and the counter are:

$$BIT == \{0, 1\}$$
$$N_{256} == 0 \mathbin{.\,.} 255$$

($N_{256}$ is just the naturals $0 \mathbin{.\,.} 2^{16} - 1$—i.e., a 16-bit unsigned.)

The computation of the counter value forms part of the predicate of schema *ComputeHitCounts*. To define a swap-out procedure, it is necessary to extend *PageFrames* a little so that information on page usage is represented.

The class and its operations are relatively straightforward.

---

**PageFrames**

$\lceil(INIT, GetPage, OverwritePhysicalPage, ClearRefBitsAndCounter,$
$\quad ComputeHitCounts, IsVictim, VictimPhysicalPageNo)$

$frames : PAGEFRAME$
$refbit : PHYSICALPAGENO \nrightarrow\mkern-14mu\rightarrow BIT$
$count : PHYSICALPAGENO \nrightarrow\mkern-14mu\rightarrow N_{256}$

$\mathrm{dom}\ count = \mathrm{dom}\ refbit$
$\mathrm{dom}\ count \subseteq \mathrm{dom}\ frames$
$\#\,\mathrm{dom}\ refbit = numrealpages$

```
┌─ INIT ─────────────────────────────────────────────────
│ (∀ i : 1 .. numrealpages •
│      frames'(i) = NullPage)
│ dom refbit' = ∅
│ dom count' = ∅
```

$GetPage \mathrel{\widehat{=}} \ldots$

$OverwritePhysicalPage \mathrel{\widehat{=}} \ldots$

$ClearRefBitsAndCounter \mathrel{\widehat{=}} \ldots$

$ComputeHitCounts \mathrel{\widehat{=}} \ldots$

$IsVictim \mathrel{\widehat{=}} \ldots$

$VictimPhysicalPageNo \mathrel{\widehat{=}} \ldots$

```
┌─ GetPage ──────────────────────────────────────────────
│ pageno? : PHYSICALPAGENO
│ fr! : PAGE
├────────────────────────────────────────────────────────
│ 1 ≤ pageno? ≤ numrealpages
│ fr! = frames(pageno?)
```

This operation retrieves a page.

The infix function *after* is required by the next schema. Its definition is repeated for convenience:

```
┌ _after_ : seq X × ℕ → seq X
├──────────────────────────────────────────────────────
│ ∀ m : seq X; offset : ℕ •
│     dom(m after offset) = (1 .. #m − offset) ∧
│     (∀ n : ℕ •
│         (n + offset) ∈ dom m ⇒ (m after offset)(n) = m(n + offset))
```

```
┌─ OverwritePhysicalPage ────────────────────────────────
│ Δ(frames)
│ pageno? : PHYSICALPAGENO
│ pg? : PAGE
├────────────────────────────────────────────────────────
│ frames' = frames ⊕ {pageno? ↦ pg?}
```

This operation overwrites a page in main store. The input *pageno?* is the index of the page frame in main store and *pg?* is a page full of data.

**Proposition 153.** *The predicate of the substitution instance of the predicate OverwritePhysicalPage[p/pageno?, pg/pg?] replaces the page indexed by p in frames by the page, pg and only that page.*

PROOF. The $\oplus$ operation can be defined as:

$$(f \oplus g)(x) = \begin{cases} f(x), x \in \operatorname{dom} f \\ g(x), \text{ otherwise.} \end{cases}$$

Then, for the predicate of the schema:

$$(frames \oplus \{pageno? \mapsto pg?\})(x) = \begin{cases} frames(x), x \in \operatorname{dom} frames \\ \{pageno? \mapsto pg?\}, x = pageno? \end{cases}$$

$\square$

The clearing of the reference bits and reference counter in a physical page is defined as follows:

```
┌─ ClearRefBitsAndCounter ──────────────────────────
│ Δ(refbit, count)
│ ppno? : PHYSICALPAGENO
├───────────────────────────────────────────────────
│ refbit' = refbit ⊕ {ppno? ↦ 0}
│ count' = count ⊕ {ppno? ↦ 0}
└───────────────────────────────────────────────────
```

The following operation computes the hit count for each page. That is, it computes the number of times the page has been referenced since it was copied into main store. It must be performed on a cyclic basis but this model does not specify how the cycle is implemented—hardware is the optimal way to compute such counts because the counter must be updated on each reference. The computation operation is defined by the following schema:

```
┌─ ComputeHitCounts ────────────────────────────────
│ Δ(pcount, count)
├───────────────────────────────────────────────────
│ (∀ i : PHYSICALPAGENO | i ∈ dom frames •
│      (∃ pcount : N₂₆₅ •
│          pcount = (count(i)/2)_mod 256 + refbit(i) * 2⁷ ∧
│          count' = count ⊕ {i ↦ pcount}))
└───────────────────────────────────────────────────
```

The lowest *count* value is chosen as the victim:

```
┌─ IsVictim ────────────────────────────────────────
│ (∃ j : PHYSICALPAGENO | j ∈ dom count •
│      count(j) = min(ran count))
└───────────────────────────────────────────────────
```

The physical page of the victim must be obtained:

```
┌─ VictimPhysicalPageNo ────────────────────────────
│ victim! : PHYSICALPAGENO
├───────────────────────────────────────────────────
│ (∃ : PHYSICALPAGENO | i ∈ dom count •
│      ∧ count(i) = min(ran count)
│      ∧ i = victim!)
└───────────────────────────────────────────────────
```

This algorithm is not foolproof but is a reasonable, hardware-independent choice. There are many alternative algorithms in the literature (see, for example, [26] or [29]) but the best will always be determined by the hardware on which the operating system runs. The assumption made here is a minimal one (because many processors implement reference bits in page frames); some machines might provide reference counters directly, while others might record the time of the last reference to each page. It is to be hoped that, in the future, victim determination will be considerably simplified by more co-operative hardware.

The *FindVictim* operation is "safe" in the sense that it will always find an in-core page. The reason for this is that the pageout operation defined above requires that $\neg$ *HaveFreePages* be true before *FindVictim* is called. This ensures that none of the candidate pages is in *freepages*.

$PGMSG ::= DELPROCPG \langle\!\langle AREF \rangle\!\rangle$
$\qquad \mid \quad GETPG \langle\!\langle AREF \times SEGMENT \times LOGICALPAGENO \rangle\!\rangle$
$\qquad \mid \quad STOPG \langle\!\langle AREF \times SEGMENT \times LOGICALPAGENO \times PAGE \rangle\!\rangle$
$\qquad \mid \quad ISPG \langle\!\langle AREF \times SEGMENT \times LOGICALPAGENO \times PAGE \rangle\!\rangle$

$FMSG ::= BADADDR \langle\!\langle APREF \times SEGMENT \times LOGICALPAGENO \rangle\!\rangle$

The ISR handling page faults can be defined as the following class. It is a subclass of the message-based ISR defined in the last chapter.

---
*PageFaultISR*

$\lceil (INIT,$
$\quad OnPageInterrupt)$

*GenericMsgISR*

---
*sched : LowLevelScheduler*
*ptab : ProcessTable*

---
*INIT*
*schd? : LowLevelScheduler*
*pt? : ProcessTable*

---
$sched' = schd? \wedge ptab' = pt?$
---

```
┌─ OnPageInterrupt ──────────────────────────────────────────────
│ intaddr? : VIRTUALADDRESS
│ ─────────────────────────────────────────────────────────────
│ ∃ fmsg : FMSG;  cp : APREF;  sg : SEGMENT;
│              lpno : LOGICALPAGENO;  offset : ℕ •
│      saddrseq(saddresstrans(intaddr?)) = sg
│   ∧ spageno(saddresstrans(intaddr?)) = lpno
│   ∧ sched.CurrentProcess[cp/cp!] ∧ sched.MakeUnready[cp/pid?]
│   ∧ ptab.DescrOfProcess[cp/pid?, pd/pd!]
│   ∧ ctxt.SwapOut ∧ fmsg = BADADDR⟪cp, sg, lpno⟫
│   ∧ (∃ isrid, fdvrid : APREF | fdrvrd = faultdrvr •
│        SendInterruptMsg[fdrvrid/driver?, fmsg/m?])
│   ⨟ ctxt.SwapIn
└────────────────────────────────────────────────────────────────
```

When a page fault occurs, the ISR sends a message to the *PageFault-Driver*. The page-fault handler or driver is defined by the following class. As is usual with classes that represent processes, it exports only one operation, here *DoOnPageFault*. The driver also contains routines that access page tables but they are not exported. The idea is that once the driver starts, it has exclusive access to these data structures. The data structures, however, still need to be protected by locks.

```
┌─ PageFaultDriver ──────────────────────────────────────────────
│ ⌈(INIT, DoOnPageFault)
│ ┌────────────────────────────────────────────────────────────
│ │ sched : LowLevelScheduler;  ptab : ProcessTable;  pts : PageTables;
│ │ vsm : VStoreManager;  pfs : PageFrames;  msgman : MsgMgr
│ │ lck : Lock
│ └────────────────────────────────────────────────────────────
│
│ ┌─ INIT ──────────────────────────────────────────────────────
│ │ mmgr? : MsgMgr;  sch? : LowLevelScheduler
│ │ ptb? : ProcessTable;  pgtabs? : PageTables
│ │ vstoreman? : VStoreManager
│ │ pgfrms? : PageFrames;  lk? : Lock
│ │ ──────────────────────────────────────────────────────────
│ │ msgman′ = mmgr? ∧ sched′ = sch? ∧ ptab′ = ptb? ∧ pts′ = pgtabs?
│ │ vsm′ = VStoreManager ∧ pfs′ = PageFrames ∧ lck′ = lk?
│ └────────────────────────────────────────────────────────────
└────────────────────────────────────────────────────────────────
```

$findVictimLogicalPage \mathrel{\widehat{=}} \ldots$

$haveVictim \mathrel{\widehat{=}} \ldots$

$findVictimPage \mathrel{\widehat{=}} \ldots$

$swapPageToDisk \mathrel{\widehat{=}} \ldots$

$retrievePageFromDisk \mathrel{\widehat{=}} \ldots$

$storePageOnDisk \mathrel{\widehat{=}} \ldots$

$genOnPageFault \mathrel{\widehat{=}} \ldots$

$onPageFault \mathrel{\widehat{=}} \ldots$

$DoOnPageFault \mathrel{\widehat{=}} \ldots$

The driver has to find a victim page to swap out when there is no free store. The following schema defines this operation. The page to be swapped out is a *logical* one at this point. The schema maps the logical page to a physical page by another operation.

___ *findVictimLogicalPage* _____

$p! : APREF$
$sg! : SEGMENT$
$lpno! : LOGICALPAGENO$
$victim! : PHYSICALPAGENO$

_____

$(\exists\, p : APREF;\ s : SEGMENT;\ l : LOGICALPAGENO\ |$
$\qquad p \in \mathrm{dom}\ pagetable \wedge sg \in \mathrm{dom}\ pagetable(p)$
$\qquad\qquad \wedge\ l \in \mathrm{dom}\ pagetable(p)(sg)\ \bullet$
$\quad pagetable(p)(sg)(l) = victim!$
$\wedge\ l \notin lockedpages(p)(sg)$
$\wedge\ p! = p \wedge sg! = s \wedge lpno! = l)$

The schema simplifies to:

_____

$p! : APREF$
$sg! : SEGMENT$
$lpno! : LOGICALPAGENO$
$victim! : PHYSICALPAGENO$

_____

$p! \in \mathrm{dom}\ pagetable \wedge sg! \in \mathrm{dom}\ pagetable(p!)$
$\wedge\ lpno! \in \mathrm{dom}\ pagetable(p!)(sg!)$
$\wedge\ pagetable(p!)(sg!)(lpno!) = victim!$
$\wedge\ lpno! = \notin lockedpages(p!)(sg!)$

The following definition is a synonym. It tests the reference count of the victim physical page.

$haveVictim \mathrel{\widehat{=}} pfs.IsVictim$

The final operation to locate a victim page is the following:

$findVictimPage \mathrel{\widehat{=}}$
    $pfs.VictimPhysicalPageNo[victim/victim!]$
        $\wedge\ findVictimLogicalPage[victim/victim!]$

It is possible to put a few properties of pages on a more formal basis.

**Proposition 154.** *Locked pages can never be victims.*

PROOF.   The predicate of *findVictimLogicalPage* contains the conjunct $l \notin$ *lockedpages*$(p)(sg)$, where $l$ is the logical page number, $p$ the process identifier and $sg$ the segment.                                                                                 □

**Proposition 155.** *Free pages can never be victims.*

PROOF.   Since free pages do not belong to any process (by Proposition 147), they do not appear in *pagetable*, so they cannot be victims because the quantifier in *findVictimLogicalPage* ranges over *pagetable*.                               □

**Proposition 156.** *Faulting processes are not ready.*

PROOF.   The predicate of *PageFaultISR.OnPageInterrupt* contains a reference to the schema *MakeUnready*$[cp/pid?]$, where $cp$ is the current process (i.e., $cp = currentp$), so it is the process that caused the page fault. The action of *MakeUnready* is to remove the process from the ready queue.                  □

**Proposition 157.** *A faulting process cannot be executed.*

PROOF.   By the previous proposition, *MakeUnready* implies that $cp$ cannot be scheduled until it is returned to the ready queue. Furthermore, $cp$ cannot continue because *OnPageInterrupt* calls *sched.ScheduleNext* to select the next process to run; this process will not be $cp$.                                                   □

**Corollary 12.** *A faulting process is blocked.*

PROOF.   This is a consequence of the previous proposition.                    □

The operation that actually swaps a page from main store to the paging disk is the following:

$swapPageToDisk \,\widehat{=}$
　　　$(findVictimPage[p/p!, sg/sg!, lpno/lpno!]$
　　　　　　$\wedge\ vsm.MarkSharedLogicalPageAsOut[p/p?, sg/sg?, lpno/lpno?]$
　　　　　　$\wedge\ pfs.GetPage[pg/fr!, victim!/pageno?]$
　　　　　　$\wedge\ storePageOnDisk[p/p?, sg/sg?, lpno/lpno?, pg/pg?]$
　　　　　　$\wedge\ pfs.ClearRefBitsAndCounter[victim!/ppno?]$
　　　　　　$\wedge\ storePageOnDisk[p/p?, sg/sg?, lpno/lpno?, pg/pg?])$
　　　　　　　　　　　　$\backslash\{pg, p, sg, lpno\}$
　　　　$\vee\ Skip$

First, a victim is located and then marked as being not in store (*MarkShared-LogicalPageAsOut*). The page causing the page fault is then demanded from the paging disk and the victim is sent to the disk for temporary storage. The reference bits of the victim are then cleared so that the referenced page (the one causing the page fault) can be written to it by *storePageOnDisk*.

　　The next operation to be defined is *retrievePageFromDisk*. This operation, as its name suggests, communicates with the paging disk process to locate the page that is to be brought into main store as a consequence of the last page fault. The reference to this page is, in fact, the one that caused the page fault that resulted in the current execution of the *PageFaultDriver*. It is defined by the following schema:

---
$\underline{\quad retrievePageFromDisk \underline{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}}$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$
$pg! : PAGE$

---
$(\exists\, msg_o, msg_r : PGMSG;\ src, dest : APREF;\ pg : PAGE\ |$
　　　　　　　　$msg_o = GETPG\langle\!\langle p?, sg?, lpno?\rangle\!\rangle$
　　　　　　　　　　$\wedge\ src = \mathsf{faultdrvr} \wedge dest = \mathsf{pagedsk}\ \bullet$
　　$msgman.SendMsg[src/src?, dest/dest?, msg_o/m?]_9^\circ$
　　$(msgman.RcvMsg[src/caller?, dest/src?, msg_r/m!]$
　　　　$\wedge\ msg_r = ISPG\langle\!\langle p?, sg?, lpno?, pg!\rangle\!\rangle))$

---

As can be seen, this operation mostly handles messages. First, the operation sends a message requesting that the paging disk retrieve the page specified by the parameters *p?*, *sg?* and *lpno?*; the page is represented by *pg!*. The page, *pg!*, is returned by the paging disk in the *ISPG* message.

　　The operation to store a page on disk is similar to the last one.

---
*storePageOnDisk*
*p*? : *APREF*
*sg*? : *SEGMENT*
*lpno*? : *LOGICALPAGENO*
*pg*? : *PAGE*

---
$(\exists\, msg : PGMSG;\ src, dest : APREF\ |$
$\qquad\qquad msg = STOPG\langle\!\langle p?, sg?, lpno?, pg?\rangle\!\rangle$
$\qquad\qquad \wedge\ src = \mathsf{faultdrvr} \wedge dest = \mathsf{pagedsk} \bullet$
$\qquad mgman.SendMsg[src/src?, dest/dest?, msg/m?])$

---

Again, this operation mostly deals with messages. In this case, it contains just one send operation. Messages of type $STOPG$ request the paging disk to store the page specified by $p?$, $sg?$, $lpno?$; the page is denoted by $pg?$.

The next operation schema defines the operation performed whenever a page fault occurs. This is a complex operation and its definition reflects this complexity. The operation receives a message containing a specification of the page that was referenced and determined (by the hardware) not to be in main store (it is on the paging disk).

The operation determines whether there are any page frames free in main store. If there are, the page is located on the disk and copied into a free page. To do this, the operation sends a message to the paging-disk process requesting the page.

If there are no page frames free in main store, a page that is currently resident in main store must be placed on the paging disk in order to make space for the one that caused the page fault. A suitable resident page must have a low frequency of reference in the near future and the reference-count mechanism specified above is used to determine which it is. This page is referred to as the *victim* in the following schema and the schemata defining the selection operations. It should be noted that the victim can belong to any process whatsoever but cannot be a page that is locked into main store. The victim is swapped to the paging disk and the one causing the page fault is retrieved and copied into the newly vacated page frame in main store. The process that caused the page fault is then returned to the ready queue and the ISR waits for the next page fault.

It should be noted that the victim can never be a locked page (i.e., a page locked into main store). This condition is imposed so that, in particular, pages locked into main store by the kernel (typically pages containing kernel code and data structures) cannot be swapped out. It is undesirable to swap kernel pages out of store because they might be involved in the operation currently being performed or they might be ISRs and the scheduler.

In some kernel designs, it is possible for kernel processes to be stored in swappable pages. In such a case, the pages would contain data or processes that are considered to be of lower importance, implementing operations that the kernel can afford to have blocked while some of their store is paged out.

The kernel assumed here is the one modelled in Chapter 4 and extended in Chapter 5. It is assumed that *all* of its components (data structures and processes) must be stored in pages that are locked into main store (and hence are stored in pages with the *locked* attribute). (The schemata defined above for victim selection, it should be noted, depend on the commutativity of conjunction.)

$genOnPageFault \; \widehat{=}$
    $\exists fmsg : FMSG; \; me, src : APREF \mid me = \mathsf{faultdrvr} \wedge src = \mathsf{hardware} \; \bullet$
        $msgman.RcvMessage[me/dest?, src/src?, fmsg/m!]$
        $\wedge \; (\exists fpid : APREF; \; fs : SEGMENT;$
                        $flpno : LOGICALPAGENO;$
                        $destpg : PHYSICALPAGENO \; \bullet$
            $fmsg = BADADDR \langle\!\langle fpid, fs, flpno \rangle\!\rangle$
            $((lck.Lock_9^\circ$
                $((pts.HaveFreePages \wedge pts.AllocateFreePage[destpg/ppno!])$
                    $\vee \; (haveVictim$
                        $\wedge \; swapPageToDisk[destpg/victim!])) \wedge lck.Unlock)_9^\circ$
            $(lck.Lock_9^\circ$
                $(retrievePageFromDisk[fpid/p?, fs/sg?, flpno/lpno?, page/pg!]$
                    $\wedge \; pfs.ClearRefBitsAndCounter[destpg/ppno?]$
                    $\wedge \; pfs.OverwritePhysicalPage[page/pg?, destpg/pageno?]$
                    $\wedge \; vsm.MarkSharedLogicalPageAsIn$
                                $[fpid/p?, fs/sg?, flpno/lpno?]$
                    $\wedge \; lck.Unlock)))$
            $_9^\circ(lck.Lock \wedge sched.MakeReady[fpid/pid?])$
            $_9^\circ(sched.ScheduleNext \wedge lck.Unlock))$


$onPageFault \; \widehat{=} \; (sched.CurrentProcess[p/cp!] \wedge genOnPageFault[p/p?]) \setminus \{p\}$


$DoOnPageFault \; \widehat{=}$
    $(\forall i : 1 \mathinner{.\,.} \infty \; \bullet \; onPageFault)$

**Proposition 158.** *If there are free pages, no page is swapped out.*

PROOF. The first component of the sequential composition in the predicate of *genOnPageFault* contains the disjunction:

$(pts.HaveFreePages \wedge pts.AllocateFreePage[destpg/ppno!])$
    $\vee \; (haveVictim \wedge swapPageToDisk[destpg/victim!])$

If there are free pages, *pts.HaveFreePages* is satisfied.                □


**Proposition 159.** *If there are no free pages in main store, a page is swapped out.*

PROOF. There is a disjunction in the composition forming *genOnPageFault*'s predicate:

$(pts.HaveFreePages \land pts.AllocateFreePage[destpg/ppno!])$
$\qquad \lor (haveVictim \land swapPageToDisk[destpg/victim!])$

In this case, if *HaveFreePages* is not satisfied, $\neg$ *HaveFreePages* must be satisfied and the swap is performed by *swapPageToDisk*. □

**Proposition 160.** *If a process fault occurs, the referenced page overwrites a page frame freed by swapping out.*

PROOF. In the previous proposition, the page *destpg* was swapped out to disk. In the second component of *genOnPageFault*'s predicate, the operation *overWritePhysicalPage* occurs as a conjunct. By $p \land q \vdash p$, the result follows. □

**Proposition 161.** *If a process faults, the referenced page is brought into store.*

PROOF. The operation *retrievePageFromDisk[page/pg]* is a conjunct in the second component of the sequential composition in *genOnPageFault*. □

**Proposition 162.** *If there are free pages, only a free page is written when a page is swapped in.*

PROOF. The predicate of schema *genOnPageFault* contains the following references:

$pts.AllocateFreePage[destpg/ppno!] \ldots {}_9^\circ \ldots$
$\qquad retrievePageFromDisk[\ldots, page/pg!]$
$\qquad pts.OverwritePhysicalPage[page/pg?, destpg/pageno?]$

The physical page, *destpg*, is the one overwritten by *page* in the operation defined by the schema *OverwritePhysicalPage*. The physical page *destpg* is allocated by *AllocateFreePage*. The page retrieved from disk is *page*; it is retrieved by *retrievePageFromDisk* (the omitted arguments refer to the faulting process). □

**Proposition 163.** *Newly swapped pages have a zero reference count.*

PROOF. Reference bits are cleared in the newly swapped page by *genOnPageFault*. The second component of this schema's predicate contains, as one of its conjuncts, a reference to *pfs.ClearRefBitsAndCounter[destpg/ppno?]*—this is a reference to a page frame in physical store, not to the contents (which is denoted by *page* in this predicate). □

**Proposition 164.** *Newly swapped pages cannot be victims unless all pages are victims.*

PROOF.  By the victim-finding operation, the victim has the minimum reference count:

```
__ IsVictim _____
 (∃ j : PHYSICALPAGENO | j ∈ dom count •
     count(j) = min(ran count))
_____
```

A newly swapped-in page—one that has not yet been referenced—has a reference count of 0. To become a victim, there must be no page with reference count $> 0$.                                                                      □

A process can be waiting on a device when one of its pages is chosen to swap out. If the driver copies data and puts it into buffers associated with the waiting process' PCBs, there is only the issue of swapping out the page. However, there is no way *a priori* of knowing whether the page just swapped out will cause a page fault when its owning process next executes. Since swapping out does not affect the operations of the device upon which the victim is waiting, it would appear valid just to pick any process that is not locked into main store.

The reader should note that, logically, this is perfectly adequate. As a proposed implementation, this is unlikely to work well. It is to be expected that page faults will be relatively frequent. The paging disk has a latency time that must be taken into account. When a user process causes a page fault, it must be blocked. Clearly, processes ought to be blocked for the shortest possible time. The paging disk, however, serialises requests. All of this suggests that the swapin/swapout operations should be as fast as possible.

Moreover, the specification, as it stands, allows the ISR to respond to as many interrupts as it can but it must also wait for the driver. The driver's message input is restricted to one immediate and one outstanding message. This suggests that the ISR should enqueue messages on the driver and immediately halt.

Furthermore, the context of the faulting process *must* be swapped immediately. This is because it cannot progress and *must* be taken off the processor: alternatives are hard to discern. Removal of the current process' registers from the processor by the ISR is, therefore, justified.

Similarly, the page disk can lose requests if it just hangs between instructing a disk search and reading the result.

This subsystem also shows limitations with the message-passing régime. Because a synchronous method has been adopted, the sender must wait if the receiver is not in a state to receive. This has the implication that, should the page disk process not be ready to accept another request (either because it is waiting for the disk or because it is processing another request), the page-fault

handler will have to wait. This has the implication that the page-fault handler might miss an interrupt.

The presence of operations to add and remove process data from the paging disk complicates matters also. Luckily, these requests will be less common than simple page faults.
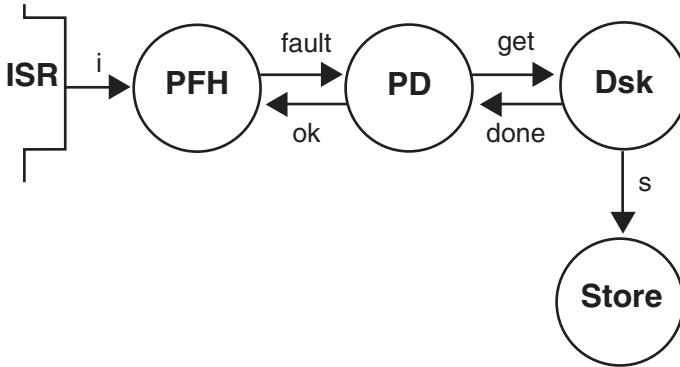


**Fig. 6.4.** *The actual specification.*

The question for us is the following. Even though the classes and operations presented above provide an adequate logical model, they do not take into consideration all of the pragmatic issues. Should the specification be altered to reflect these pragmatic issues?

The specification presented above can be represented diagrammatically as in Figure 6.4. As can be seen, the virtual-storage management subsystem consists of the ISR (denoted by *ISR* in the figure), the page-fault driver process (denoted by *PFH* in the figure) and the disk driver process (denoted by *PD*). The figure also includes the disk itself (denoted by *DSK*) and main store (denoted by *STORE*).

The arrows in Figure 6.4 denote the messages or interactions between processes. The arrow labelled *i* denotes the message sent by the *ISR* to the page-fault handling process (*PFH*). This message contains the specification of the page that was referenced (in terms of the segment and logical page numbers) and the process (its *APREF*). The page-fault handler sends a *fault* message containing the same information to the paging-disk process, which in turn sends a request (as a *get* message) to the disk proper (actually, to the disk driver). The disk driver retrieves the page and sends a *done* message to the paging-disk handler. The *done* message denotes the fact that the retrieval has been successfully completed.

Once the page has been written to main store, process *PD* sends an *ok* message to the page-fault handler process, *PFH*. On reception of the *ok*, the page-fault handler can wait for another page fault.

It is clearly a highly desirable property for any optimisations of the basic (logical) specification to behave in the same way as the specification. It is also highly desirable, given the present context, to be able to demonstrate this in a formal way. In order to achieve this, the proposed optimisations are formalised as CCS [21] processes so that they can be manipulated in formally sound ways. CCS is chosen as the representation because we are interested in the interactions between the component processes of the subsystem, not in the specification of the components. The processes to be modelled do not have properties suggestive of the use of the $\pi$-calculus (e.g., mobility), so CCS appears sufficient.

The subsystem can be represented in CCS as the following set of equations:

$ISR = \bar{i}.ISR$
$PFH = i.\overline{fault}.ok.PFH$
$PD = fault.\overline{get}.done.\bar{s}.\overline{ok}.PD$
$DSK = get.\overline{done}.DSK$
$STORE = s.STORE$

The overall arrangement is represented in CCS as:

$VM_1 = (ISR \mid PFH \mid PD \mid DSK \mid STORE) \setminus \{i, fault, ok, get, done, s\}$

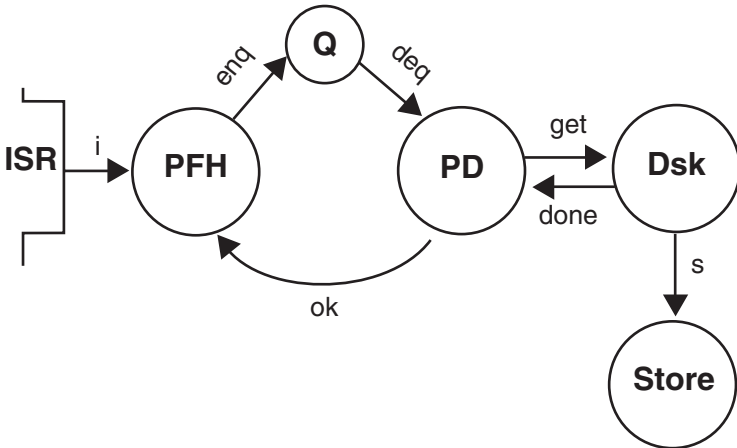(Note that actions are hidden using the $\setminus$ operation.)



**Fig. 6.5.** *The specification using a queue.*

As noted above, the design depicted in Figure 6.4 and represented by the above set of CCS process definitions can be optimised. An obvious optimisa-

tion is to introduce a queue of requests between the page-fault handler ($PFH$) and paging-disk process ($PD$). This is the arrangement shown in Figure 6.5. In the arrangement shown in Figure 6.5, the $PFH$ process places new requests into the queue. The queue is represented by the process named $Q$ in the figure, and the operation of sending a request (really, just an enqueuing of the request) is represented by the *enq* arrow. Requests are removed from the queue process by the paging-disk process, $PD$, in exactly the same way they are in the first case. When the page has been copied to main store, the paging-disk process, $PD$ sends the page-fault handling process an *ok* to inform it that: (i) the copy has been performed and that (ii) the process that caused the fault just rectified can now be unblocked.

The argument is that this second version can process more page faults per unit time. In this case, $PFH$ does not now wait for the page fault to be rectified before it can wait for a new fault. Instead, it passes the request to the rest of the subsystem and then immediately blocks on the page-fault ISR.

This second arrangement can be represented by CCS processes as follows:

$$VM_2 = init.(ISR \mid PFH_2 \mid Q \mid PD_2 \mid DSK \mid STORE)$$
$$\setminus \{i, enq, deq, ok, done, get, init, s\}$$

For the definition of this subsystem, processes $ISR$, $DSK$ and $STORE$ remain as in the first case. The remaining processes must be redefined as follows (the subscripts will be explained below).

$PFH_2 = i.\overline{enq}.ok.PFH_2$
$Q = init.enq.Q_1$
$Q_1 = enq.Q1 + deq.Q1$
$PD_2 = deq.\overline{get}.done.\overline{s}.\overline{ok}.PD_2$

It is clearly necessary to distinguish between the two versions of $PFH$ that have been defined at this point (a third version will be added shortly). For this reason, subscripts were introduced into the specifications.

It would be useful for these two specifications (models) to be equivalent in some sense. One important sense is that they should be *observationally* equivalent; another is that they should be *bisimilar*.

The property of *observational* equivalence of two processes is very much the intuitive one: two processes are observationally equivalent when they cannot be distinguished by an external observer. In other words, the externally visible events that can be perceived by an external observer are determined by the observer to be the same in content and in order, no matter which process is observed. In the cases of $VM_1$ and $VM_2$, the externally observable events are restricted by the hiding operator ($\setminus$, as in Z).

Bisimilarity is an equivalence that also takes hidden actions into account. (Those readers unfamiliar with the concept should consult [21], Chapter 4, for an extended treatment.)

It is now possible to engage in formal reasoning about processes $VM_1$ and $VM_2$ and to prove two important propositions about them. Rather than

engaging in a hand proof, it was considered interesting to employ automation. To this end, the *Concurrency Workbench of the New Century* [8] was employed as a tool. The Concurrency Workbench can be used to determine a number of properties of CCS and CSP processes, including observational equivalence and bisimilarity. The propositions concerning equivalence of the various versions of the subsystem were all proved using the Concurrency Workbench.

**Proposition 165.** *Processes $VM_1$ and $VM_2$ are observationally equivalent.*

PROOF. Both $VM_1$ and $VM_2$ were encoded in the input format required by the CWB and tested using the `eq -S obseq` command. The CWB system determined that $VM_1$ and $VM_2$ are observationally equivalent.          □

Matters are different when it comes to bisimilarity. Processes $VM_1$ and $VM_2$ are quite dissimilar in internal structure, and process $VM_2$ offers an initial event, *init*, which initialises the process $Q$, so it does not appear, at first sight, that $VM_1$ and $VM_2$ will be bisimilar. Indeed, this is the case.

**Proposition 166.** *Processes $VM_1$ and $VM_2$ are not bisimilar.*

PROOF. Both $VM_1$ and $VM_2$ were encoded in the input format required by the CWB and tested using the `eq -S bisim` command. The CWB system determined that $VM_1$ and $VM_2$ are *not* bisimilar.          □
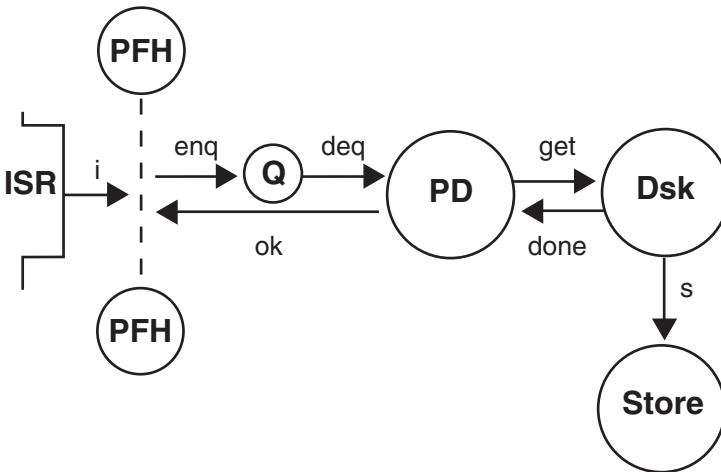


**Fig. 6.6.** *The ideal specification.*

Finally, a second specification can be considered as an optimisation. This version retains the queue of requests. It differs from the second version by creating a new instance of the *PFH* process whenever an interrupt occurs. This instance performs the same operations as in the second version but immediately terminates (denoted by the **0** process). The CCS specification is (again, subscripts are used to differentiate between versions):

$$VM_3 = init.(ISR \mid PFH_3 \mid Q \mid PD_2 \mid DSK \mid STORE)$$
$$\backslash\{i, enq, deq, ok, done, get, init, s\}$$

where:

$$PFH_3 = i.(\overline{enq}.ok.\mathbf{0} \mid PFH_3)$$

**Proposition 167.** *$VM_3$ is observationally equivalent to $VM_2$.*

PROOF. The models for $VM_2$ and $VM_3$ were encoded in the input format required by the CWB and tested using the `eq -S obseq` command. The CWB system determined that the two are observationally equivalent.     □

**Proposition 168.** *$VM_3$ is bisimilar to $VM_2$.*

PROOF. The models for $VM_2$ and $VM_3$ were encoded in the input format required by the CWB and tested using the `eq -S obseq` command. The CWB system determined that the two are bisimilar.     □

From the first of these two results, we have the next proposition:

**Proposition 169.** *Process $VM_1$ is observationally equivalent to $VM_3$.*

PROOF. By transitivity of the equivalence relation.

(In addition, the encoded models were tested using the `eq -S obseq` command and the result was supported by the CWB system.)     □

We also have the following negative result:

**Proposition 170.** *Processes $VM_1$ and $VM_3$ are not bisimilar.*

PROOF. Again, by transitivity of the equivalence relation, this result would be expected.

(In addition, the CWB supported the claim.)     □

The last result is not as bad as it sounds. All that is required, here, is that the processes *appear* the same as far as an external observer is concerned. This property is sufficient for the optimisations to be considered equivalent to the original. However, as far as this chapter is concerned, the first version (the one referred to as the "logical" one) is the one that will be adopted.

With this discussion of optimisation out of the way, it is possible to return to the main theme. There follow some propositions dealing with the properties of the logical (Z) model of page faults.

**Proposition 171.** *Page faults are serviced in the order in which they occur.*

PROOF.   There are two senses:

1. Each page fault is an interrupt. In this sense, page faults are all processed in order.
2. In this sense, page faults are serviced by the handler process which services messages sent to it by the page-fault ISR. These messages are sent in temporal (sequence) order. That the handler process operates upon these messages in the correct order is a consequence of the correctness of the message-passing operations.

Each sense implies the statement of the proposition.                    □

**Proposition 172.** *Every page fault is serviced eventually (i.e., the page-fault mechanism is fair).*

PROOF.   This proof follows from the correctness and fairness of the message-passing operations.                    □

**Proposition 173.** *The process causing a page fault has status* pstrunning *when the fault occurs.*

PROOF.   Only a running process can cause a page fault.

ISRs cannot cause faults and are not processes.

Drivers are processes. They are locked into main store and, by hypothesis, have all the necessary pages locked in as well.

Only user processes have pages that can be swapped into and out of store.

The fact that a fault occurs implies that the process causing it is executing. There can only be one process that is executing at any one time (i.e., is the current value of *currentp* and has a status of pstrunnning).

This is a property of the scheduler.                    □

**Proposition 174.** *A faulting process is not unblocked until the replacement page has been swapped into physical store.*

PROOF.   The service operation *onPageFault* first blocks the faulting process by setting its state to pstwaiting and by removing it from the processor. Therefore *blockFaultingProcess* implies that $currentp \neq currentp' \land regs' \neq regs$, where *regs* are hardware registers.

The *onPageFault* operation takes the form of a sequential composition. There are three operations, thus forming a composition of the form $S_1 \mathbin{\text{\textsubscript{9}}} S_2 \mathbin{\text{\textsubscript{9}}} S_3$. The second operation is itself a sequential composition; it is this part that performs the swapping operation. The operations are, therefore, totally ordered and the order can be directly related to temporal succession.

Let $p$ denote the faulting process and let the composition be written as above. Then post $S_1 \Rightarrow status(p) = \mathsf{pstwaiting}$ and $currentp' \neq p$.

Throughout $S_2$, $status(p) = \mathsf{pstwaiting}$.

Finally, pre $S_3 \Rightarrow status(p) = \mathsf{pstwaiting}$, while post $S_3 \Rightarrow status(p) = \mathsf{pstready}$ but $p \neq currentp$—the last being a property of *MakeReady*. It cannot be true because the service process is currently executing. □

**Proposition 175.** *If a process, p, causes a page fault, that process is not marked ready until the faulting page has been replaced.*

PROOF.  This follows immediately from the previous proposition.           □

**Proposition 176.** *After execution of onPageFault, there is exactly one physical page in physical store such that its logical page number in the faulting process maps to the physical page number.*

PROOF.  The proposition implies that exactly one page is mapped into store by the page-fault mechanism.

Inspection of the predicate of *genOnPageFault* shows that only one page is introduced into store.           □

**Proposition 177.** *If a process, p, causes a page fault, after that page fault has been serviced, the process is in the ready state.*

PROOF.  That the process is blocked follows immediately from Proposition 174. The predicate of *genOnPageFault* contains an instance of *MakeReady* with the faulting process' identifier substituted for the formal parameter. □

**Proposition 178.** *If a page fault occurs and a page has to be swapped out, that page can be a page belonging to the faulting process or to another process.*

PROOF.  There are two cases to consider.
Case 1: If there are free pages, they are consumed. Only the faulting process is affected.
Case 2: *HaveVictim* is called. The victim page is found in physical store by *pfs.IsVictim*:

$findVictimPage \mathrel{\widehat{=}}$
    $pfs.VictimPhysicalPageNo$
    $\wedge\ findVictimLogicalPage$

The value returned by *pfs.VictimPhysicalPageNo* is just a physical page number with the minimum reference count. The schema does not mention processes. Similarly, *findVictimLogicalPage* in this class merely looks up the

owning process' identifier and the segment in which the page occurs, as well as the logical page number of the victim page. Therefore, the page can belong to *any* process in the system except those that lock it into main store (i.e., driver and system processes).                                              □

**Proposition 179.** *If a page fault occurs and there are free pages in physical store, physical store is updated; only one process is affected.*

PROOF.  By the definition of *onPageFault*:

$$pts.HaveFreePages \wedge pts.AllocateFreePages$$

implies that the physical page to which the swapped-in page is to be written is a free page in physical store.

In this case, only the faulting process is affected. This is because the allocation of a free page relates only to a single process (the process to which the allocation is made).                                              □

**Corollary 13.** *If a page fault occurs within process, p, and if a page has to be swapped out, only a maximum of two processes are affected by the page-swapping operation.*

PROOF.  Immediate from the two preceding propositions.                □

**Proposition 180.** *No pages are swapped unless a page fault occurs.*

PROOF.  Inspection reveals that only the schema *genOnPageFault* references the relevant operations.                                              □

### 6.3.4 Extending Process Storage

Processes very often make requests to increase the amount of store they use. In a paged system, this allocates more pages to the process. This subsection contains a model of the operations. It is rarer for processes to release store but operations to perform this task are defined as well. Pages can also be shared between processes. Sharing can be used, for example, to implement message passing in virtual storage.

First, there is a problem with allocating pages. The problem is that there might not be a physical page free when the request is made. The system could block the requesting process until there is a free page in the page frame. This solution does not fit well with the aims of virtual storage. Furthermore, when a new process is created, it will request a set of pages to hold code, data, stack,

etc. It would not be acceptable to block the creation of the process until main store becomes free.

A solution more in keeping with the purpose of virtual store is the following. When a request is made to allocate a new page and there is no free physical store, an empty page is allocated on disk. The new page is allocated to the requesting process and can be written to by the requesting process.

This mechanism can be used when allocating extra pages to a new page.

Therefore, it is necessary to begin with the definition of the empty page:

$$
\begin{array}{l}
pageofzeroes : PAGE \\
\hline
\forall\, pg : PAGE \bullet \\
\quad pageofzeroes = \lambda\, i : 1 \mathbin{..} framesize \bullet 0
\end{array}
$$

This structure need not be held in a page; it can be produced in a loop and then set into a buffer of smaller size. Ideally, it should be packed into a single disk block and then handed to the paging disk. This needs to be iterated $n$ times, where $n$ is the number of disk blocks per page.

Pages must be specified when allocating, deallocating and sharing. A method for specifying them is required. Here, we define a structure for this purpose:

$$
PAGESPEC == APREF \times SEGMENT \times LOGICALPAGENO
$$

The following axiomatic definitions are required to manipulate objects of type *PAGESPEC*. The relation that comes first will be discussed below. The rest of the definitions are: the constructor function (*mkpgspec*) and the accessor functions (*pgspecpref* to obtain the process reference, *pgspecseg* to obtain the segment name and *pgspeclpno* to obtain the logical page number).

$$
\begin{array}{l}
smap : PAGESPEC \leftrightarrow PAGESPEC \\
mkpgspec : APREF \times SEGMENT \times LOGICALPAGENO \rightarrow PAGESPEC \\
pgspecpref : PAGESPEC \rightarrow APREF \\
pgspecseg : PAGESPEC \rightarrow SEGMENT \\
pgspeclpno : PAGESPEC \rightarrow LOGICALPAGENO \\
\hline
\forall\, p : APREF;\ sg : SEGMENT;\ lpno : LOGICALPAGENO \bullet \\
\quad mkpgspec(p, sg, lpno) = (p, sg, lpno) \\
\\
\forall\, ps : PAGESPEC \bullet \\
\quad pgspecpref(ps) = fst\ ps \\
\quad pgspecseg(ps) = fst(snd\ ps) \\
\quad pgspeclpno(ps) = snd^{2}\ ps
\end{array}
$$

The type *PAGESPEC* could have been defined earlier in this chapter. If it had been used there, it would have been necessary to represent and manipulate all virtual addresses and page references in terms of *PAGESPEC*. Although this seems attractive, we believe, it would have complicated the specifications somewhat.

The axiomatic definitions begin with *smap*, a relationship between elements of *PAGESPEC*. This is the *sharing map*. When two elements of *PAGESPEC* are in the relation, the processes mentioned as the first component of *PAGESPEC* share the page referred to by the two *PAGESPEC*s.

The operations to allocate, deallocate and share pages are collected into a (somewhat *ad hoc*) class called *VStoreManager*. The class is intended to act as a component in an interface library. It exports most of its operations so that a wide variety of combined operations can be defined elsewhere.

---
**VStoreManager**

$\lceil(INIT, AddNewMainStorePageToProcess, AddNewVirtualPageToProcess,$
$\quad CanAddPageToProcess, MarkLogicalPageAsShared, UnshareLogicalPage,$
$\quad WithdrawLogicalPage, SharedLogicalPageSharers, IsSharedLogicalPage,$
$\quad RemoveSharedLogicalPageOwner, RemoveLogicalPageSharer,$
$\quad RawShareLogicalPageBetweenProcesses, ShareLogicalPageBetweenProcesses,$
$\quad ReturnSharedLogicalPageToOwner, MarkSharedLogicalPageAsIn,$
$\quad MarkSharedLogicalPageAsOut, ShareLogicalSegment,$
$\quad ReleaseSharedSegment, ReleaseSegmentPagesExcept,$
$\quad CanReleaseSegment, SharedPagesInSegment, CanReleaseProcessVStore)$

---
$pts : PageTables$
$pfs : PageFrames$

---
**INIT**

$pgtabs? : PageTables$
$pfrms? : PageFrames$

---
$pts' = pgtabs?$
$pfs' = pfrms?$

---

$makeEmptyPage \mathrel{\widehat{=}} \ldots$

$AddNewMainStorePageToProcess \mathrel{\widehat{=}} \ldots$

$AddNewVirtualPageToProcess \mathrel{\widehat{=}} \ldots$

$CanAddPageToProcess \mathrel{\widehat{=}} \ldots$

$MarkLogicalPageAsShared \mathrel{\widehat{=}} \ldots$

$UnshareLogicalPage \mathrel{\widehat{=}} \ldots$

$WithdrawLogicalPage \mathrel{\widehat{=}} \ldots$

$SharedLogicalPageSharers \mathrel{\widehat{=}} \ldots$

$IsSharedLogicalPage \mathrel{\widehat{=}} \ldots$

$RemoveSharedLogicalPageOwner \mathrel{\widehat{=}} \ldots$

$RemoveLogicalPageSharer \mathrel{\widehat{=}} \ldots$

$RawShareLogicalPageBetweenProcesses \mathrel{\widehat{=}} \ldots$

$ShareLogicalPageBetweenProcesses \mathrel{\widehat{=}} \ldots$

$ReturnSharedLogicalPageToOwner \mathrel{\widehat{=}} \ldots$

$MarkSharedLogicalPageAsIn \mathrel{\widehat{=}} \ldots$

$MarkSharedLogicalPageAsOut \mathrel{\widehat{=}} \ldots$

$ShareLogicalSegment \mathrel{\widehat{=}} \ldots$

$ReleaseSharedSegment \mathrel{\widehat{=}} \ldots$

$ReleaseSegmentPagesExcept \mathrel{\widehat{=}} \ldots$

$CanReleaseSegment \mathrel{\widehat{=}} \ldots$

$SharedPagesInSegment \mathrel{\widehat{=}} \ldots$

$CanReleaseProcessVStore \mathrel{\widehat{=}} \ldots$

First, there is the hidden operation that creates an empty page:

---
$makeEmptyPage$

$pg! : PAGE$

---
$pg! = pageofzeroes$

---

This operation is hidden because it is somewhat undesirable for everyone to manipulate the buffers in which new pages are created.

The operation to add a new page to a process using a free-store page is the following. It allocates the page frame, adds the page to the process and clears the page (this is not really necessary but is a nice feature[2]).

$AddNewMainStorePageToProcess \mathrel{\widehat{=}}$
$\quad (pts.HaveFreePages$
$\qquad \wedge ((pts.IncProcessPageCount \mathbin{\substack{\circ \\ 9}} (pts.LatestPageCount)$
$\qquad\qquad \wedge pts.AllocateFreePage[ppno/ppno!]$
$\qquad \wedge pts.AddPageToProcess[ppno/ppno?, lpno!/lpno?]$
$\qquad \wedge makeEmptyPage[pg/pg?]$
$\qquad \wedge pts.OverwritePhysicalPage[ppno/pageno?, pg/pg?]) \setminus \{ppno\}))$

Note that an alternative is to swap out a process' page. The approach here is simpler but much more profligate.

The next operation uses the prceding one and then stores the page that it has created on the paging disk. Note that the operation requires there to be *at least* one page-sized buffer in the kernel.

$AddNewVirtualPageToProcess \mathrel{\widehat{=}}$
$\quad (\neg pts.HaveFreePages$
$\qquad \wedge (pts.IncProcessPageCount \mathbin{\substack{\circ \\ 9}} pts.LatestPageCount)$

---

[2] Some (civilised?) operating systems perform this operation on allocating new store.

$$\wedge\,(\exists\,ppno : PHYSICALPAGENO \mid ppno = 1 \bullet$$
$$(makeEmptyPage[pg/pg!]$$
$$\wedge\,pts.AddPageToProcess[lpno!/lpno?, ppno/ppno?]$$
$$\wedge\,StorePageOnDisk[lpno!/lpno?, pg/pg?])) \setminus \{pg\})$$

The value assigned to *ppno* is purely arbitrary. The physical page number of any page on disk has no relevance because it will be mapped to another physical page when swapped into main store.

The following is a predicate. It is true if the process denoted by *p*? can add at least one page to its *sg*? segment.

---
*CanAddPageToSegment* _____
*p*? : *APREF*
*sg*? : *SEGMENT*

$pagecount(p?)(sg?) < maxvirtpagespersegment$

---

The actual operation to add a new page to a process is the following. Depending upon the state of main store, it adds the page directly or stores it on the paging disk:

$AddNewPageToProcess \,\widehat{=}$
    $AddNewMainStorePageToProcess$
        $\vee\,AddNewVirtualPageToProcess$

**Proposition 181.** *If there are no free pages in main store, a zero page is created for it and written to disk.*

PROOF.  The schema *AddNewPageToProcess* is a disjunction:

$AddNewMainStorePageToProcess$
    $\vee\,AddNewVirtualPageToProcess$

Each disjunct is a conjunction, with mutually exclusive first conjuncts.

In particular, *AddNewVirtualPageToProcess* has a predicate containing the following:

$makeEmptyPage[pg/pg!]$
    $\wedge$
    $\ldots$
    $\wedge\,StorePageOnDisk[\ldots, pg/pg?]$

The existential quantifier can be removed by the one-point rule, and the result follows from the fact that $p \wedge q \Rightarrow p$.                    □

**Proposition 182.** *If there are free pages, one is allocated for the newly created page.*

PROOF. Similar to the above, concentrating on *AddNewMainStorePage*. Again, $p \wedge q \vdash p$ allows the conclusion to be drawn.    □

The only thing that needs to be done when a page is shared between processes is to mark it. This is actually done by adding *PAGESPEC*s to the sharing map, *smap*:

---
*MarkLogicalPageAsShared*

$\Delta(smap)$
$ownproc?, shareproc? : APREF$
$ownseg?, shareseg? : SEGMENT$
$ownlp?, sharelp? : LOGICALPAGENO$

---

$(\exists\, atrip1, atrip2 : PAGESPEC \,\bullet$
$\qquad atrip1 = mkpgspec(ownproc?, ownseg?, ownlp?) \wedge$
$\qquad atrip2 = mkpgspec(shareproc?, shareseg?, sharelp?) \wedge$
$\qquad smap' = smap \cup \{(atrip1, atrip2)\})$

---

Unsharing, conversely, removes a pair of *PAGESPEC*s from the sharing map:

---
*UnshareLogicalPage*

$\Delta(smap)$
$p? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

---

$smap' = smap \vartriangleright \{mkpgspec(p?, sg?, lpno?)\}$

---

If a page is to be reallocated or removed from store completely, it must be totally removed from the sharing map:

---
*WithdrawLogicalPage*

$\Delta(smap)$
$owner? : APREF$
$sg? : SEGMENT$
$lpno? : LOGICALPAGENO$

---

$smap' = \{mkpgspec(p?, sg?, lpno?)\} \vartriangleleft smap$

---

In later operations, it is necessary to determine which processes share a given page. The page is specified by the *sg?* and *lpno?* inputs and is shared by (at least) *p?*. The relational image is used to determine all the sharing *PAGESPEC*s and, hence, the sharing processes:

```
┌─ SharedLogicalPageSharers ──────────────────────────────
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
│ srs! : 𝔽 PAGESPEC
├─────────────────────────────────────────────────────────
│ smap⦇ {mkpgspec(p?, sg?, lpno?)} ⦈ = srs!
└─────────────────────────────────────────────────────────
```

The following predicate is true when the page is shared:

```
┌─ IsSharedLogicalPage ───────────────────────────────────
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
├─────────────────────────────────────────────────────────
│ mkpgspec(p?, sg?, lpno?) ∈ dom smap
└─────────────────────────────────────────────────────────
```

When, say, a process terminates or when a shared page used to contain a message is withdrawn, the entry for the page must also be removed from the *smap*:

```
┌─ RemoveSharedLogicalPageOwner ──────────────────────────
│ Δ(smap)
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
├─────────────────────────────────────────────────────────
│ smap' = {mkpgspec(p?, sg?, lpno?)} ⩤ smap
└─────────────────────────────────────────────────────────
```

All the sharers of a page can be removed from the *smap* by the following operation:

```
┌─ RemoveLogicalPageSharers ──────────────────────────────
│ Δ(smap)
│ p? : APREF
│ sg? : SEGMENT
│ lpno? : LOGICALPAGENO
├─────────────────────────────────────────────────────────
│ (∃ pg : PAGESPEC | pg = mkpgspec(p?, sg?, lpno?) •
│      smap' = smap ▷ {pg})
└─────────────────────────────────────────────────────────
```

The following pair of schemata define the operation of sharing a page between two processes. The first schema defines the actual operation, while the second hides the logical page number. The first operation returns the logical page number of the shared page. Sometimes, this is something of a nuisance, so the second operation hides it. The first operation returns the page number so that it can be referenced, while the second hides the page

number—approximately, it states that a page is shared but does not say which it is.

$RawShareLogicalPageBetweenProcesses \; \widehat{=}$
    $((pts.IsLockedPage[owner?/p?, ownseg?/sg?, ownlpno?/lpno?]$
        $\wedge \; pts.LockPage[sharer?/p?, shareseg?/sg?, sharelpno!/lpno?])$
        $\vee \; Skip)$
    $\wedge \; NextNewLogicalPageNo[sharer?/p?, shareseg?/sg?, sharelpno!/lpno!]$
    $\wedge \; (pts.PhysicalPageNo[owner?/p?, ownseg?/sg?, ownlpno?/lpno?, ppno/ppno!]$
    $\wedge \; pts.AddPageToProcess[sharer?/p?, shareseg?/sg?, sharelpno!/lpno?]$
    $\wedge \; MarkLogicalPageAsShared[owner?/ownproc?, sharer?/shareproc?,$
                        $ownlpno?/ownlp?, sharelpno!/sharelp?]$
    $\wedge \; pts.MarkPageAsShared[owner?/p?, ownseg?/sg?, ownlpno?/lpno?]$
    $\wedge \; pts.MarkPageAsShared[sharer?/p?, shareseg?/sg?, sharelpno!/lpno?])$
                $\backslash \{ppno\}$


$ShareLogicalPageBetweenProcesses \; \widehat{=}$
    $RawShareLogicalPageBetweenProcesses \setminus \{sharelpno!\}$

**Proposition 183.** *If two processes share a page, pg, that page will appear in the page tables of both processes.*

PROOF.   The page, *pg*?, is already assumed to be in the page table of its owning process. The predicate of *RawShareLogicalPageBetweenProcesses* contains, *AddPageToProcess*[*sharer?/p?, shareseg?/sg?, sharelpno!/lpno?*] as a conjunct. The arguments to the substitution refer to the process that is receiving the page.  □


**Proposition 184.** *Schema RawShareLogicalPageBetweenProcesses makes the specified page shared by both processes.*

PROOF.  The predicate of *RawShareLogicalPageBetweenProcesses* contains an instance of *MarkPageAsShared*, which simplifies to $smap' = smap \cup \{(s_1, s_2)\}$, where:
$$s_1 = mkpgspec(ownproc?, ownseg?, ownlp?)$$
and
$$s_2 = mkpgspec(sharer?, shareseg?, sharelp?)$$
and where *ownproc*? is the owning process identifier and *sharer*? is the identifier of the process being granted the right to share the page; *ownseg*? and *shareseg*? denote the segments in the owner and sharer's spaces; *ownlp*? is the owning process' logical page number and *sharelp*? is the logical page number in the sharer's segment.  □

The only point to note about the following is that it unshares the page in question.

$ReturnSharedLogicalPageToOwner \;\widehat{=}$
      $pts.IsSharedPage[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$
      $\land\; pts.UnsharePage[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$
      $\land\; pts.RemovePageProperties[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$
      $\land\; pts.RemovePageFromPageTable[sharer?/p?, shareseg?/sg/,$
                                 $sharelpno?/lpno?]$
      $\land\; UnshareLogicalPage[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$

**Proposition 185.** *When a page is unshared, it is removed from one of the page tables.*

PROOF.  The operation *ReturnSharedLogicalPageToOwner* removes the specified page from one of the processes that share it. This is done by removing the page from the sharer's page table by means of the conjunct:

$$RemovePageFromPageTable[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$$

The other conjuncts cancel various attributes, in particular:

$$RemovePageProperties[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$$

as well as removing the information that this page is shared by process *sharer?* by
$$UnsharePage[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$$
and

$$UnshareLogicalPage[sharer?/p?, shareseg?/sg/, sharelpno?/lpno?]$$

. $\hfill\square$

Shared pages need to be marked as in-store or out-of-store. The following pair of schemata define these operations. They will be used below.

$MarkSharedLogicalPageAsIn \;\widehat{=}$
      $(pts.IsSharedPage \land pts.MarkPageAsIn \land$
         $(SharedLogicalPageSharers[srs/srs!] \land$
            $(\forall pg : PAGESPEC \mid pg \in srs \bullet$
               $(\exists p : APREF;\; sg : SEGMENT;$
                        $lpno : LOGICALPAGENO \bullet$
             $pg = mkpgspec(p, sg, lpno)$
              $\land\; pts.MarkPageAsIn[p/p?, sg/sg?, lpno/lpno?])) \setminus \{srs\}))$
      $\lor\; pts.MarkPageAsIn$

**Proposition 186.** *When a shared page is swapped in, it is swapped in in all sharing processes.*

PROOF. The operation *MarkPageAsIn* occurs within the scope of the universal quantifier which ranges over elements of *PAGESPEC*. These elements are exactly those describing the shared page in processes other than its owner. The operation *MarkPageAsIn* is applied to each descriptor (rather, to its components) and records the fact that the page has entered main store in each of the processes that share it.                                                                    □

$$
\begin{aligned}
&MarkSharedLogicalPageAsOut \mathrel{\widehat{=}} \\
&\quad (pts.IsSharedPage \wedge pts.MarkPageAsOut \wedge \\
&\qquad (SharedLogicalPageSharers[srs/srs!] \wedge \\
&\qquad (\forall\, pg : PAGESPEC \mid pg \in srs \bullet \\
&\qquad\quad (\exists\, p : APREF;\ sg : SEGMENT; \\
&\qquad\qquad\qquad\qquad lpno : LOGICALPAGENO \bullet \\
&\qquad\qquad mkpgspec(p, sg, lpno) = pg \wedge \\
&\qquad\qquad pts.MarkPageAsOut[p/p?, sg/sg?, lpno/lpno?]))) \setminus \{srs\}) \\
&\quad \vee\ pts.MarkPageAsOut
\end{aligned}
$$

**Proposition 187.** *When a shared page is swapped out, it is swapped out in all sharing processes.*

PROOF. The operation defined by *MarkPageAsOut* occurs inside the universal quantifier in the predicate of *MarkSharedLogicalPageAsOut*. The universal quantifier ranges over all *PAGESPEC* elements that are related to the process in which the page is being swapped. Each element of type *PAGESPEC* records the process identifier, segment identifier and logical page number of the page in the sharing processes. Therefore, the quantifier ranges over all descriptions of the page in the sharing process. The *MarkPageAsOut* operation is applied to each of these descriptions, thus recording the fact that the page is swapped out.                                                                    □

**Proposition 188.** *The operation MarkSharedLogicalPageAsOut can never mark a locked page as swapped out.*

PROOF. By Proposition 154.                                                       □

The next schema defines the operation to return the next logical page number available in the segment, *sg*?, of process *p*? (both *sg*? and *p*? are inputs to the schema):

$$
\begin{aligned}
&NextNewLogicalPageNo \mathrel{\widehat{=}} \\
&\quad pts.IncProcessPageCount \mathbin{\S} pts.LatestPageCount
\end{aligned}
$$

Segments, too, can be shared and unshared. This is not obviously useful but it turns out to be. When a process spawns a copy of itself or creates a child process that shares its code, the code segment can be shared if it is read-only (executable is a page attribute defined at the start of this chapter).

---
__*ShareLogicalSegment*_____

$ownerp?, sharerp? : APREF$
$ownerseg?, sharerseg? : SEGMENT$

---

$(\forall\, lpno : LOGICALPAGENO \mid$
$\qquad\qquad lpno \in \operatorname{dom} pagetable(ownerp?)(sharerp?) \bullet$
$\qquad (\exists\, ps_1, ps_2 : PAGESPEC;\; nxtlpno : LOGICALPAGENO \bullet$
$\qquad\qquad ps_1 = mkpgspec(ownerp?, ownerseg?, lpno) \wedge$
$\qquad\qquad NextNewLogicalPageNo[nxtlpno/lpno!] \wedge$
$\qquad\qquad\quad ps_2 = mkpgspec(sharerp?, sharerseg?, nxtlpno) \wedge$
$\qquad\qquad pts.MarkPageAsShared$
$\qquad\qquad\qquad [ownerp?/p?, ownerseg?/sg?, lpno/lpno?] \wedge$
$\qquad\qquad MarkLogicalPageAsShared[lpno/ownlp?, nxtlpno/sharelp?]))$

---

---
__*ReleaseSharedSegment*_____

$\Delta(pagetable)$
$p? : APREF$
$sg? : SEGMENT$

---

$pagetable' = \{sg?\} \lhd pagetable(p?)$

---

---
__*ReleaseSegmentPagesExcept*_____

$\Delta(pagetable)$
$p? : APREF$
$sg? : SEGMENT$
$except? : \mathbb{F}\ LOGICALPAGENO$

---

$pagetable' = ((\operatorname{dom} pagetable(p?)(sg?)) \setminus except?) \lhd pagetable(p?)(sg?)$

---

**Proposition 189.** *If except? = $\varnothing$, ReleaseSegmentPagesExcept implies that the predicate of CanReleaseSegment is satisfied.*

PROOF.  The predicate of *CanReleaseSegment* is:

$$\operatorname{dom} pagetable(p?)(sg?) = \varnothing$$

The predicate of *ReleaseSegmentPagesExcept* is:

$$pagetable' = ((\operatorname{dom} pagetable(p?)(sg?)) \setminus except?) \lhd pagetable(p?)(sg?)$$

Clearly (substituting $\varnothing$ for *except?*):

$$((\operatorname{dom} pagetable(p?)(sg?)) \setminus \varnothing) = \operatorname{dom} pagetable(p?)(sg?)$$

Now,

$$(\mathrm{dom}\,pagetable(p?)(sg?)) \lhd pagetable(p?)(sg?)$$

implies that dom $pagetable(p?)(sg?) = \varnothing$. This is the predicate of the schema *CanReleaseSegment*, so we are done.                                              □

**Proposition 190.** *A process whose entire virtual store has been released has no pages.*

Proof.   Obvious given the last proposition and the appropriate schema definitions.                                                                                     □

The following predicate is true if and only if all the pages in the segment, $sg?$, have been deallocated:

```
__ CanReleaseSegment _____
 p? : APREF
 sg? : SEGMENT
_____
 dom pagetable(p?)(sg?) = ∅
```

The pages in the segment $sg?$ of process $p?$ are output by this operation. They constitute the set $sps!$.

```
__ SharedPagesInSegment _____
 p? : APREF
 sg? : SEGMENT
 sps! : 𝔽 LOGICALPAGENO
_____
 (∃ psg : 𝔽 PAGESPEC •
     psg = smap⦇ {pg : PAGESPEC; lpn : LOGICALPAGENO |
                  pg = mkpgspec(p?, sg?, lpn) • pg} ⦈ ∧
     sps! = {pg : PAGESPEC; lpno : LOGICALPAGENO | pg ∈ pgs •
              pgspeclpno(pg)})
```

If a process has no pages and no segments (i.e., they have never been allocated or all have been released), its page-table entry can be released. The following schema is a predicate defining this case.

```
__ CanReleaseProcessVStore _____
 p? : APREF
_____
 dom pagetable(p?) = ∅
```

# 6.4 Using Virtual Storage

## 6.4.1 Introduction

This section is a rather looser model than the above model. It deals with the user-level interfaces. The view of virtual store is that much more akin to real

store: a linear sequence of locations that can be addressed sequentially. The concept of the virtual address is also different in this section. Instead of a complex structure, virtual addresses are considered atomic, in effect a subset of $\mathbb{N}$.

### 6.4.2 Virtual Addresses

From this point on, the specification will deal with virtual addresses, not with logical and physical pages. This amounts to the user process' perspective on virtual storage. Hitherto, all the perspective has been that of the virtual storage constructor, so virtual address spaces and virtual addresses have been seen in terms of their components.

For this reason, the following is required:

$$AllocatePageReturningStartVAddress \mathrel{\widehat{=}}$$
$$AddNewPageToProcess \wedge ComputePageVAddress[lpno!/lpno?]$$

This operation, as its name suggests, allocates a page and returns the virtual address of its start. (Note that the logical page number is *not* hidden; there are reasons for this, as will be seen.)

The following operation computes a virtual address from its components.

$$
\boxed{
\begin{array}{l}
\underline{\ ComputePageVAddress\ }\\
lpno? : LOGICALPAGENO\\
vaddr! : \mathbb{N}\\
\hline
vaddr! = (lpno? - pgallocstart) * framesize
\end{array}
}
$$

Here, *pgallocstart* is a constant. It is defined as:

$$
\boxed{
\begin{array}{l}
pgallocstart : \mathbb{N}\\
\hline
pgallocstart = 0
\end{array}
}
$$

The value of 0 is completely arbitrary, as is now explained.

Some systems map a virtual copy of the operating system onto the virtual address space of each user space (and some privileged processes, too). The virtual copy can be allocated at the start or at the end of virtual store. The allocation of user-process pages has to respect this. Furthermore, there are special addresses that are reserved by the hardware; for example, interrupt vectors, device buffers and status words. These, too, must be allocated somewhere that cannot be directly accessed by untrusted user processes.

The *pgallocstrt* represents the logical page number of the first page in a segment that can be allocated.

For clarity and simplicity, it will be assumed here that the operating system virtual copy is located entirely in its own space. Transfer to the operating

system from user space is achieved by a "mode switch". The mode switch activates additional instructions, for example those manipulating interrupts and the translation lookaside buffer. How the mode switch is performed is outside the scope of this book for reasons already given. Mode switches are, though, very common on hardware that supports virtual store.

In some systems, there are parts of the kernel space that are shared between all processes in the system. These pages are pre-allocated and added to the process image when it is created. Because they are pre-allocated, the allocation of user pages in that process must be allocated at some page whose logical page number is greater than zero. The constant *pgallocstart* denotes this offset.

Usually, the offset is used in the data segment only. For simplicity, the offset is here set to 0. Moreover, it is uniformly applied to all segments (since it is 0, this does not hurt).

The one hard constraint on virtual store is that some *physical* pages must never be allocated to user space. These are the pages that hold the device registers and other special addresses just mentioned.

Virtual-store pages are frequently marked as:

- execute only (which implies read-only);
- read-only;
- read-write.

Sometimes, pages are marked write-only. This is unusual for user pages but could be common if device buffers are mapped to virtual store pages.

The operations required to mark pages alter the attributes defined at the start of this chapter. The operations are relatively simple to define and are also intuitively clear. They are operations belonging to the class defined below in Section 6.5.2; in the meanwhile, they are presented without comment.

$MarkPageAsReadOnly \cong$
    $MakePageReadable \wedge$
    $((IsPageWritable \wedge MakePageNotWritable) \vee$
    $((IsPageExecutable \wedge MakePageNotExecutable)))$

$MarkPageAsReadWrite \cong$
    $(MakePageReadable \fatsemi MakePageWritable) \wedge$
    $(IsPageExecutable \wedge MakePageNotExecutable)$

$MarkPageAsCode \cong$
    $(MakePageExecutable \fatsemi MakePageReadable) \wedge$
    $(IsPageExecutable \wedge MakePageNotExecutable)$

An extremely useful, but generic, operation is the following. It allocates $n$ pages at the same time:

$AllocateNPages \,\widehat{=}$
  $(\forall\, p : 1 \ldots numpages? \bullet$
    $(p = 1 \wedge AllocatePageReturningStartVAddress[startaddr?/vaddr!])$
    $\vee\, (p > 1 \,\wedge$
      $(\exists\, vaddr : VADDR \bullet$
        $AllocatePageReturningStartVAddress[vaddr/vaddr!])))$

As it stands, this operation is not of much use in this model. The reason for
this is that it does not set page attributes, in particular the read-only, execute-
only and the read-write attributes. For this reason the following operations are
defined. The collection starts with the operation for allocating $n$ executable
pages.

The following allocation operation is used when code is marked as exe-
cutable and read-only. This is how Unix and a number of other systems treat
code.

$AllocateNExecutablePages \,\widehat{=}$
  $(\forall\, p : 1 \ldots numpages? \bullet$
    $(p = 1 \wedge AllocatePageReturningStartVAddress[startaddr?/vaddr!]$
      $\wedge MarkPageAsCode)$
    $\vee\, (p > 1 \,\wedge$
      $(\exists\, vaddr : VADDR \bullet$
        $AllocatePageReturningStartVAddress[vaddr/vaddr!] \wedge$
        $MarkPageAsCode)))$

Similarly, the following operations allocate $n$ pages for the requesting pro-
cess. It should be remembered that the pages might be allocated on the paging
disk and not in main store.

$AllocateNReadWritePages \,\widehat{=}$
  $(\forall\, p : 1 \ldots numpages? \bullet$
    $(p = 1 \wedge AllocatePageReturningStartVAddress[startaddr?/vaddr!]$
      $\wedge MarkPageAsReadWrite)$
    $\vee\, (p > 1 \,\wedge$
      $(\exists\, vaddr : VADDR \bullet$
        $AllocatePageReturningStartVAddress[vaddr/vaddr!] \wedge$
        $MarkPageAsReadWrite)))$

$AllocateNReadOnlyPages \,\widehat{=}$
  $(\forall\, p : 1 \ldots numpages? \bullet$
    $(p = 1 \wedge AllocatePageReturningStartVAddress[startaddr?/vaddr!]$
      $\wedge MarkPageAsReadOnly)$
    $\vee\, (p > 1 \,\wedge$
      $(\exists\, vaddr : VADDR \bullet$
        $AllocatePageReturningStartVAddress[vaddr/vaddr!] \wedge$
        $MarkPageAsReadOnly)))$

To support the illusion of virtual storage, virtual addresses can be thought of as just natural numbers (including 0):

$$VADDR == \mathbb{N}$$

and the virtual store for each process can be considered a potentially infinite sequence of equal-sized locations:

```
┌─ VirtualStore ──────────────────────────────────
│  vlocs : seq PSU
│  maxaddr : ℕ
│ ─────────────────────
│  #locs = maxaddr
└──────────────────────────────────────────────────
```

It must be emphasised that each virtual address space has *its own copy* of the *VirtualStore* schema. Page and segment sharing, of course, make regions of store belonging to one virtual address space appear (by magic?) as part of another.

The usual operations (read and write) will be supported. However, when the relevant address is not present in real store, a *page fault* occurs and the *OnPageFault* driver is invoked with the address to bring the required page into store.

For much of the remainder, a block-copy operation is required. This is used to copy data into pages based on addresses. For the time being, it can be assumed that every address is valid (the hardware should trap this, in any case).

```
┌─ CopyVStoreBlock ───────────────────────────────
│  ΔVirtualStore
│  data? : seq PSU
│  numelts? : ℕ
│  destaddr? : VADDR
│ ─────────────────────
│  vlocs' = (λ i : 1 .. (destaddr? − 1) • vlocs(i))
│              ⌢ ⟨data?⟩ ⌢ (vlocs after (destaddr? + numelts?))
└──────────────────────────────────────────────────
```

If any of the addresses used in this schema are not in main store, the page-faulting mechanism will ensure that it is loaded.

Operations defining the user's view of virtual store are collected into the following class. It is defined just to collect the operations in one place. (In the next section, the operations are not so collected—they are just assumed to be part of a library and are, therefore, defined in Z.)

The class is defined as follows. The definition is somewhat sparse and contains only two operations, *CopyVStoreBlock* and *CopyVStoreFromVStore*. In a full implementation, this class could be extended considerably. The point, here, though, is merely to indicate that operations similar to those often implemented for real store can be implemented in virtual storage systems at

a level above that at which virtual addresses are manipulated as complex
entities.

```
┌─ Users VStore ──────────────────────────────────────────────
│ ⌈(INIT, CopyVStoreBlock, CopyVStoreFromVStore)
│ ┌──────────────────────────────────────────────────────────
│ │ vlocs : seq PSU
│ │ maxaddr : ℕ
│ │ ─────────────────
│ │ #locs = maxaddr
│ └──────────────────────────────────────────────────────────
│
│ ┌─ INIT ───────────────────────────────────────────────────
│ │ maxaddress? : ℕ
│ │ ──────────────────────────
│ │ maxaddr′ = maxaddress?
│ └──────────────────────────────────────────────────────────
│
│ CopyVStoreBlock ≘ ...
│ CopyVStoreFromVStore ≘ ...
└──────────────────────────────────────────────────────────────
```

Data can be copied into a virtual-store page (by a user process) by the
following operation:

```
┌─ CopyVStoreBlock ───────────────────────────────────────────
│ ΔVirtualStore
│ data? : seq PSU
│ numelts? : ℕ
│ destaddr? : VADDR
│ ─────────────────────────────────────────────────────────────
│ vlocs′ = (λ i : 1 .. (destaddr? − 1) • vlocs(i))
│                  ⌢⟨data?⟩ ⌢ (vlocs after (destaddr? + numelts?))
└──────────────────────────────────────────────────────────────
```

A similar operation is the following. It copies one piece of virtual store
to another. It is useful when using pages as inter-process messages: the data
comprising the message's payload can be copied into the destination (which
might be a shared page in the case of a message) from the page in which it
was assembled by this operation:

```
┌─ CopyVStoreFromVStore ──────────────────────────────────────
│ ΔVirtualStore
│ fromaddr? : VADDR
│ toaddr? : VADDR
│ numunits? : ℕ₁
│ ─────────────────────────────────────────────────────────────
│ (∃ endaddr : VADDR | endaddr = fromaddr? + numunits? − 1 •
│     vlocs′ = (λ i : 1 .. (toaddr? − 1) • vlocs(i))
│                  ⌢(λ j : fromaddr? .. endaddr • vlocs(j))
│                  ⌢(vlocs after endaddr + 1))
└──────────────────────────────────────────────────────────────
```

### 6.4.3 Mapping Pages to Disk (and Vice Versa)

Linux contains an operation called `memmap` in its library. This maps virtual store to disk store and is rather useful (it could be used to implement persistent store as well as other things, heaps for instance).

A class is defined to collect the operations together. Again, this class is intended only as an indication of what is possible. In a real system, it could be extended considerably; for example, permitting the controlled mapping of pages between processes, archiving of pages, and so on.

$$\boxed{\begin{array}{l} \underline{\;PageMapping\;} \\ \lceil(INIT, MapPageToDisk, MapPageFromDiskExtendingStore) \\[4pt] \hline \begin{array}{l} usrvm : UserVStore \\ pfr : PageFrames \end{array} \\[6pt] \hline \begin{array}{l} \underline{\;INIT\;} \\ uvm? : UserVStore \\ pgfrm? : PageFrames \\ \hline usrvm' = uvm? \\ pfr' = pgfrm? \end{array} \\[6pt] MapPageToDisk \,\widehat{=}\, \ldots \\ MapPageFromDiskExtendingStore \,\widehat{=}\, \ldots \\ writePageToDisk \,\widehat{=}\, \ldots \\ readPageFromDisk \,\widehat{=}\, \ldots \end{array}}$$

The operations in this class are all fairly obvious, as is their operation. Commentary is, therefore, omitted.

$$\boxed{\begin{array}{l} \underline{\;writePageToDisk\;} \\ \ldots \\ diskparams? : \ldots \\ pg? : PAGE \\ \hline \ldots \end{array}}$$

$$
\begin{aligned}
MapPageToDisk \,\widehat{=}\, & \\
& (pgt.PhysicalPageNo[ppno/ppgno!] \,\wedge \\
& \quad pfr.GetPage[ppno/pageno?, pg/fr!] \,\wedge \\
& \quad writePageToDisk[pg/pg?]) \setminus \{ppno, pg\}
\end{aligned}
$$

```
┌─ readPageFromDisk ────────────────────────────────────────────
│ . . .
│ diskparams? : . . .
│ pg! : PAGE
│ ────────────────────
│
│ . . .
└───────────────────────────────────────────────────────────────
```

This operation extends the store of the requesting process. It is used when reading pages from disk. The disk page is added to the process' virtual storage image.

$MapPageFromDiskExtendingStore \mathrel{\hat{=}}$
$\qquad usrvm.AllocatePageReturningStartVAddress[pageaddr!/vaddr!]\,_9^\circ$
$\qquad\qquad (readPageFromDisk[pg/pg!] \wedge$
$\qquad\qquad (\exists\, sz : \mathbb{N} \mid sz = framesize \bullet$
$\qquad\qquad\qquad usrvm.CopyVStoreBlock$
$\qquad\qquad\qquad\qquad\qquad [sz/numelts?, pg/data?, pageaddr!/destaddr?]))$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \backslash\{pg\}$

There ought to be an operation to delete the page from the image. However, only the most careful programmers will ever use it, so the operation is omitted. It is, in any case, fairly easy to define.

Note that there is no operation to map a disk page onto an existing virtual-store page. This is because it will probably be used extremely rarely.

The operations in this class could be extended so that the specified disk as well as the paging disk get updated when the frame's counter is incremented. This would automatically extend the disk image. A justification for this is that it implements a way of protecting executing processes from hardware and software failure. It can be used as a form of journalling.

This scheme can also be used on disk files. More generally, it can also work on arbitrary devices. This could be an interesting mechanism to explore when considering virtual machines of greater scope (it is an idea suggested by VME/B). Since this is just speculation, no more will be said on it.

### 6.4.4 New (User) Process Allocation and Deallocation

This section deals *only* with user-process allocation and deallocation. The general principles are the same for system processes but the details might differ slighty (in particular, the default marking of pages as read-only, etc.).

When a new process is created, the following schema is used. In addition, the virtual-store-management pages must be set up for the process. This will be added to the following schema in a compound definition.

─────────────────────────────────────────

**UserStoreMgr**

⌈(*INIT*, *MarkPageAsReadOnly*, *MarkPageAsReadWrite*, *MarkPageAsCode*,
   *AllocateNPages*, *AllocateNExecutablePages*, *AllocateNReadWritePages*,
   *AllocateNReadOnlyPages*, *CopyVStoreBlock*, *CopyVStoreFromVStore*,
   *AllocateNewProcessStorage*, *ReleaseSharedPages*,
   *FinalizeProcessPages*, *AllocateCloneProcessStorage*)

─────────────────────────────────────────

*usrvm* : *UserVStore*
*pgt* : *PageTables*

─────────────────────────────────────────

  **INIT**

  *uvm?* : *UserVStore*
  *ptbl?* : *PageTables*
  ──────────────
  $usrvm' = uvm?$
  $pgt' = ptbl?$

─────────────────────────────────────────

$MarkPageAsReadOnly \,\widehat{=}\, \ldots$

$MarkPageAsReadWrite \,\widehat{=}\, \ldots$

$MarkPageAsCode \,\widehat{=}\, \ldots$

$AllocateNPages \,\widehat{=}\, \ldots$

$AllocateNExecutablePages \,\widehat{=}\, \ldots$

$AllocateNReadWritePages \,\widehat{=}\, \ldots$

$AllocateNReadOnlyPages \,\widehat{=}\, \ldots$

$CopyVStoreBlock \,\widehat{=}\, \ldots$

$CopyVStoreFromVStore \,\widehat{=}\, \ldots$

$AllocateNewProcessStorage \,\widehat{=}\, \ldots$

$ReleaseSharedPages \,\widehat{=}\, \ldots$

$FinalizeProcessPages \,\widehat{=}\, \ldots$

$AllocateCloneProcessStorage \,\widehat{=}\, \ldots$

─────────────────────────────────────────

**AllocateNewProcessStorage**

*p?* : *APREF*
*codepages?* : seq *PSU*
*codesz?*, *stacksz?*, *datasz?*, *heapsz?* : $\mathbb{N}$

─────────────────────────────────────────

$(\exists\, sg : SEGMENT;\; codeszunits : \mathbb{N}\mid$
               $sg = \mathsf{code} \land codeszunits = \#codepages? \bullet$

$AllocateNExecutablePages$
$\qquad\qquad [codesz?/numpages?, addr/startaddr!]^9_9$
$\qquad usrvm.CopyVStoreBlock$
$\qquad\qquad [codepages?/data?, codeszunits/numelts?,$
$\qquad\qquad\qquad addr/destaddr?])$
$(\exists\, sg : SEGMENT \mid sg = \mathsf{data} \bullet$
$\quad AllocateNReadOnlyPages[datasz?/numpages?])$
$(\exists\, sg : SEGMENT \mid sg = \mathsf{stack} \bullet$
$\quad AllocateNReadWritePages[stacksz?/numpages?])$
$(\exists\, sg : SEGMENT \mid sg = \mathsf{heap} \bullet$
$\quad AllocateNReadWritePages[heapsz?/numpages?])$

---

$\underline{\quad ReleaseSharedPages\ }$
$\Delta(smap)$
$p? : APREF$

$(\forall\, ps : PAGESPEC \mid ps \in \mathrm{dom}\, smap \wedge pgspecpref(ps) = p? \bullet$
$\quad (\exists\, s : PAGESPEC \mid (ps, s) \in smap \bullet$
$\qquad smap' = smap \setminus \{(ps, s)\}))$
$\wedge\, (\forall\, ps : PAGESPEC \mid ps \in \mathrm{ran}\, smap \wedge psgpecpref(ps) = p? \bullet$
$\quad (\exists\, s : PAGESPEC \mid (s, ps) \in smap \bullet$
$\qquad smap' = smap \setminus \{(s, ps)\}))$

---

Once this schema has been executed, the process can release all of its pages:

$FinalizeProcessPages \;\widehat{=}$
$\qquad pgt.RemovePageProperties \wedge pgt.RemoveProcessFromPageTable$

---

$\underline{\quad AllocateCloneProcessStorage\ }$
$p? : APREF$
$clonedfrom? : APREF$
$stacksz?, datasz?, heapsz? : \mathbb{N}$

$(\exists\, sg : SEGMENT \mid sg = \mathsf{code} \bullet$
$\quad ShareLogicalSegment[clonedfrom?/ownerp?, p?/sharerp?,$
$\qquad\qquad\qquad sg/ownerseg?, sg/sharerseg?])$
$(\exists\, sg : SEGMENT \mid sg = \mathsf{data} \bullet$
$\quad AllocateNReadOnlyPages[datasz?/numpages?])$
$(\exists\, sg : SEGMENT \mid sg = \mathsf{stack} \bullet$
$\quad AllocateNReadWritePages[stacksz?/numpages?])$
$(\exists\, sg : SEGMENT \mid sg = \mathsf{heap} \bullet$
$\quad AllocateNReadWritePages[heapsz?/numpages?])$

This works because of the following argument. The first of the two schemata (*ReleaseSharedPages*) above first removes the process from all of the pages that it shares but does not own. Then it removes itself from all of those shared pages that it does own. This leaves it with only those pages that belong to it and are not shared with any other process.

If a child process performs the first operation, it will remove itself from all of the pages it shares with its parents; it will also delete all of the pages it owns. The parent is still in possession of the formerly shared pages, which might be shared with other processes. As long as the parent is blocked until all of its children have terminated, it cannot delete a page that at least one of its children uses. Thus, when all of a process' children have terminated, the parent can terminate, too. Termination involves execution of the operations defined by *ReleaseSharedPages* and *FinalizeProcessPages*.

The only problem comes with clones. If the clone terminates before the original, all is well. Should the original terminate, it will delete pages still in use by the clone. Therefore, the original must also wait for the clone to terminate.

An alternative—one that is possible—is for the owner to "give" its shared pages to the clone. Typically, the clone will only require the code segment and have an empty code segment of its own. If the code segment can be handed over to the clone in one operation (or an *atomic* operation), the original can terminate without waiting for the clone or clones. Either is possible.

The allocation of child processes is exactly the same as cloning. The difference is in the treatment of the process: is it marked as a child or as a completely independent process? Depending upon the details of the process model, a child process might share code with its parent (as it does in Unix systems), whereas an independent process will tend to have its own code (or maybe a copy of its creator's code). In all cases, the data segment of the new process, as well as its stack, will be allocated in a newly allocated set of pages. In this chapter's model, data and stack will be allocated in newly allocated segments. The mechanisms for sharing segments of all kinds have been modelled in this chapter, as have those for the allocation of new segments (and pages). The storage model presented in this chapter can, therefore, support many different process models.

## 6.5 Real and Virtual Devices

There is often confusion between real and virtual devices. It is sometimes thought that the use of virtual store implies the use of virtual devices. This is not so. In most operating systems with virtual store, the devices remain real, while in some real-store operating systems, devices are virtual.

Virtual devices are really interfaces to actual, real ones. Virtual devices can be allocated on the basis of one virtual device to each process. The virtual device sends messages to and receives them from the device process. Messages

are used to implement requests and replies in the obvious fashion. Messages to the real device from the virtual devices are just enqueued by the device process and serviced in some order (say, FIFO).

The interface to the virtual device can also abstract further from the real device. This is because virtual devices are just pieces of software. For example, a virtual disk could just define read and write operations, together with return codes denoting the success of the operation. Underneath this simple interface, the virtual device can implement more complex interfaces, thus absolving the higher levels of software from the need to deal with them. This comes at the cost of inflexibility.

This model can be implemented quite easily using the operations already defined in this book. Using message passing, it can be quite nicely structured.

There is another sense in which devices can be virtualised. Each device interface consists of one or more addresses. Physical device interfaces also include interrupts. Operations performed on these addresses control the device and exchange data between device and software. The addresses at which the device interface is located are invariably fixed in the address map. However, in a virtual system, there is the opportunity to map the pages containing device interfaces are mapped into the address space of each process. (This can be done, of course, using the sharing mechanism defined in this chapter.) This allows processes directly to address devices. However, some form of synchronisation must be included so that the devices are fairly shared between processes (or virtual address spaces). Such synchronisation would have to be included within the software interface to each device and this software can be at as low a level as desired.

A higher-level approach is to map standard addresses (by sharing pages) into each address space but to include a more easily programmed interface. Again, the mechanisms defined in this book can be used as the basis for this scheme.

## 6.6 Message Passing in Virtual Store

At a number of points in this chapter, the idea of using shared pages (or sets of shared pages) to pass messages between processes has been raised. The basic mechanisms for implementing message passing have also been defined.

When one process needs to send a message to another, it will allocate a page and mark it as shared with the other process. Data will typically be placed in the page before sharing has been performed. The data copy operation can be performed by one of the block-copy operations, *CopyVStoreBlock* or *CopyVStoreFromVStore* (Section 6.5.1).

The receiving process must be notified of the existence of the new page in its address space. This can be achieved as either a synchronous or an asynchronous event—the storage model is completely neutral with respect to

this. In a system with virtual storage, message passing will be implemented as system calls, so notification can be handled by kernel operations. For example, the synchronous message-passing primitives defined in Chapter 5 can easily be modified to do this. What is required is that the message call point to a page and not to a small block of storage. Equally, the asynchronous mechanism outlined in Chapter 3 can be modified in a similar fashion.

Message passing based on shared pages will be somewhat slower at runtime than a scheme based upon passing pointers to shared storage blocks (buffers), even when copying buffers between processes is required. The reason for this is clear from an inspection of the virtual storage mechanisms. For this reason, it would probably be best to implement two message-passing schemes: one for kernel and one for user messages. The kernel message scheme would be based on shared buffers within kernel space; user messages would use the shared-page mechanism outlined above.

In some cases, additional system processes are required in addition to those executing inside the kernel address space and they will be allocated their own virtual store. In order to optimise message passing between these processes and the kernel, a set of pages can be declared as shared but not *incore* (i.e., not locked into main store). The set of pages can be pre-allocated by the kernel at initialisation time, so no new pages need to be allocated. All that remains is for the pages to be given to the processes. This can be achieved using the primitives defined in this chapter.

## 6.7  Process Creation and Termination; Swapping

Process creation, activation and termination are unaffected by the virtual storage mechanisms. The virtual storage subsystem must be booted before any processes are created, so all processes, even those inside the kernel, are created in virtual address space. The primitives to allocate and deallocate storage have been defined above (Sections 6.5.1 and 6.5.3). The operations to create and delete processes can be implemented in a way analogous to those defined in Chapter 4 (and assumed in Chapter 5), with the virtual-store primitives replacing those handling real store. The most significant difference between the two schemes is that the virtual-store allocation operations are not as limited in the amount of store they can allocate. The virtual storage operations are only limited by the number of pages permitted in a segment and not be the size of main store.

Virtual store also has advantages where swapping is concerned. It is possible to include a swapping system in a virtual-store-based system. As with the scheme defined in detail in Chapter 4, the swapper will transfer entire process images between main store and the swap disk (or swap file). Under virtual storage, the swapper treats the page as the basic unit for transfer. The swapper reads the page table and swaps physical pages to disk. Not all segments need be swapped to disk; code segments might be retained in main store while

there are active child processes. The process is, however, complicated by the fact that a process image is likely to be shared between the paging disk and main store.