

作業系統設計與實務

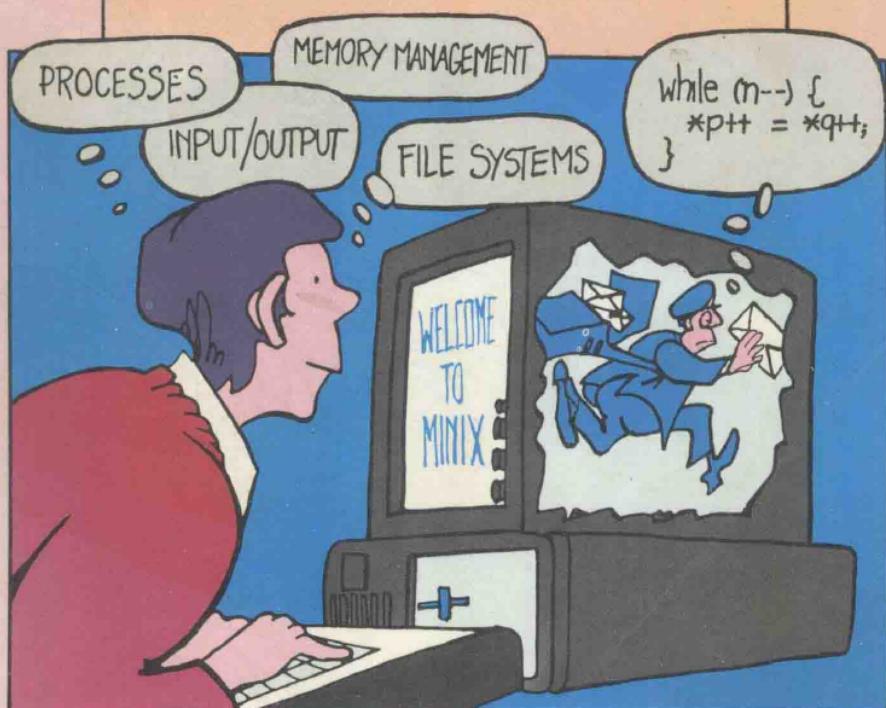
譯 者

胡 紹 杰 莊 達 三

OPERATING SYSTEMS

DESIGN AND IMPLEMENTATION

Andrew S. Tanenbaum



東華書局印行

作業系統

設計與實務

原著者

Andrew S. Tanenbaum

譯 者

胡紹杰 · 莊達三

東華書局印行



版權所有・翻印必究

經原出版公司授權獨家翻譯

中華民國七十六年十二月初版

作業系統設計與實務

定價 新臺幣叁佰伍拾元整

著者 Andrew S. Tanenbaum

譯者 胡紹杰・莊達三

發行人 卓 鑑 蘇

出版者 臺灣東華書局股份有限公司
臺北市博愛路一〇五號
郵撥：00064813

印刷者 合興印刷廠

行政院新聞局登記證 局版臺業字第零柒貳伍號
(76074)

OPERATING SYSTEMS:

Design
and
Implementation

原序

大部分關於作業系統的書籍，都是重於理論而缺乏實際應用；然而，本書在這兩方面都很平衡，它不但包含了所有基礎理論的細節，像是行程、行程與行程間的通訊、信號機（semaphore）、監督程式（monitor）、訊息分頁法（message paging）、遠端程序呼叫、排班的演繹法、輸入和輸出部分、死結（deadlocks）、裝置驅動程式（device driver）、記憶體管理法、分頁法、檔案系統的設計、網路上的檔案服務程式（network file servers）、自動處理（auomatic transaction）、安全和保護的技巧等等；它同時也討論了 MINIX，（一個和 UNIX 相容的作業系統），不但講解得十分詳細，並且提供了完整的原始程式供讀者研習。這樣的安排，不僅使讀者學習了理論部分，同時更可以了解到這些理論是如何被實際應用在真正的作業系統上。

一個作業系統有四個主要部分：包括行程安排(process management)，輸入和輸出，記憶體管理，和檔案系統。本書對上述的四個主題，各有一章做專門的介紹，而每一章中對所有相關的理論也都有相當的探討，並從不同系統中舉例說明，這些系統包括了 UNIX，MS-DOS，CP/M，MUTICS，和其他的作業系統。我在作業系統上有十五年的教學經驗，並曾擔任三種不同作業系統工程師（這三種機器分別為 PDP-11, 68000, IBM PC），這使我了解到如果只專門探討理論（死結、排班法則等等），會使學生對主題有相當地曲解。大部分書籍花了很多時間在排班法則（scheduling algorithm）的探討上，但事實上如果忽略了 I/O 部分（這部分通常佔有系統的百分之三十或更多），那麼排班法就可能僅是一頁的程式而已。

爲了要改正這種不平衡，我寫了一個新的作業系統，就是 MINIX。MINIX 有和第七版 UNIX 完全相同的系統呼叫（除了刪除數量很小的不重要部分）。此外，我也提供了和 UNIX 的 shell 功能上完全相同的 shell，及六十多種與 UNIX 相類似的其它程式（例如：cat，cc，cp，grep，ls，和 make）。簡而言之，對使用者而言，MINIX 和 UNIX 是非常的相像；但就內容而言，這系統完全是新的。我花了相當大的心血來架構這個系統，使它易於了解，並且讓學生可以簡單地修改它。這系統的主要部分是由分離的幾個模組（module）構成，它們彼此以訊息傳遞的方式來聯絡。所有的程序（procedure）都不長，並且都是結構化的。在原始程式中有超過三千個的註解。就像所有的作業系統一樣，MINIX 分成四個主要部分：行程管理（process management），輸入和輸出（裝置的驅動程式）（device drivers），記憶體管理（memory management）和檔案系統。這四個主題在本書中各佔了一章。每一章的開始先探討一般性的理論，然後再告訴你這主題在 MINIX 上是如何處理完成的。

完整的 MINIX 原始程式有幾種形式，也可以用不同的方式來使用。第一種方式用於 IBM PC，XT 或者 AT（或是與它們相容的產品），它將二元碼和原始程式放在磁碟片上。如果有足夠的 PC 數，那麼每個學生都可有一份這軟體的拷貝，並在他自己的 PC 上進行修改和測試。這種方式下，最好能有 640K 的 RAM 和兩個 360K 的軟性磁碟機，如果是比其要小的環境，也是可以的。硬式磁碟機並不需要。另一種方式是給每一個學生一份在 VAX 機器或其他分時系統電腦上 MINIX 的原始程式拷貝，當然學生可以在上述之一的機器上進行修改和編譯的工作，以產生可以在 IBM PC 上執行的二元機器碼，最後再將此機器碼轉錄到 IBM PC 上。這樣的安排，主要是因爲學生多，而 IBM PC 的數目卻不夠。這方式需要在學生使用的主機器上，有一個能產生 IBM PC 機器碼的 C 語言編譯器。像這樣的編譯器，在 Amsterdam Compiler Kit 上有介紹（參見 Tanenbaum, et al 所寫的文章，刊登在 Communications of the ACM Sept. 1983, pp. 654-660）。讀者不妨參考一下。另外，MINIX 在沒有 IBM PC 的情況下，也可以當做課程來使用。我們有兩種可能的

方式：第一種是由 Prentice-Hall 提供的 IBM PC 模擬程式（儲存於磁帶中），這模擬程式一次解譯一個 8088 的指令，檢查並執行它們，就像一個真實的 8088 在工作一樣。它同時也模擬了 MINIX 上的 I/O 設備。另一種可能是使學生能只用檔案系統來工作。MINIX 的檔案系統事實上只是一個大的 C 語言程式，這程式像用戶程式一樣，在作業系統外執行。它是一個和 UNIX 相容的遠端檔案服務程式，藉由訊息傳遞方式和在相同機器上的 MINIX 其它部分聯絡（我們可以很容易的設定它，使它像一個真正的網路檔案服務程式一般）。

為了使學生在主要電腦（host computer）上用這套檔案系統來實驗，檔案系統和學生的測試程式應該分開編譯。一個主程式需要建立管道、分岔，和執行測試程式與檔案系統，它們經由管道來傳遞訊息與聯絡。而在一個執行 UNIX 的 VAX 系統上，做這些事所需要的軟體都在 Prentice Hall 磁帶上。由於這檔案系統，建立程式，和 IBM PC 模擬都是用 C 寫的，所以可以很容易的把它們放到非 VAX 系統上。至於那些不能對 MINIX 進行修改或實驗的情況，學生仍然可以經由研讀原始程式（見本書附錄）來學習 MINIX。

最後，對那些只在理論上有興趣而在 MINIX 上卻沒有興趣的讀者而言，討論 MINIX 執行的章節（已很清楚地作上記號），可以跳過不看也不會接不下去。在這種情形下，這本書也可以作為一本平常用的作業系統教科書，不必參考 MINIX。至於那些有 USENET 權利，並希望投遞更多的軟體，建議改進等等的讀者，也可以在 comp.os.minix (MINIX 作業系統公司) 上做。而在課程使用上，一本解答手冊是需要的。我們可以從 Prentice-Hall 訂購。

在這計劃的進行過程期間，我非常幸運地得到許多人的幫忙。我要謝謝 Dick Grune, Wiebren de Jonge, Jan Looyen, Jim van Keulen, Hans van Staveren, Jennifer Steiner, 和 Peter Weinberger，他們研讀部分的原稿並給我們許多建議。也特別感謝 Brian Kernighan 研讀 1.4 及經常提醒我關於法則 13 的各種細節。雖然我寫了這作業系統原始程式的 12000 行，但也有一些由其它人提供的公用程式。沒有他們的幫助，MINIX 不可能有用。另外，我要謝謝 Martin Atkins, Erik Baalbergen,

charles Forsyth , Richard Gregg , Paul Polderman , 和 Robbert van Renesee 。 Paul Ogilvie 將這系統放到 MS-DOS 之下。 Michiel Huisjes , Patrick van Kleef , 和 Adri Koppes 則提供大量的幫助在這軟體的許多地方，這些都遠遠超過了他們的職責，但他們還是做了。最後，我要謝謝 Suzanne ，因為她永無止盡的耐心，沒有她的支持和了解，我沒有辦法在 PC 面前完成 MINIX 。我也要謝謝 Barbara 和 Marvin 用 Suzanne 的電腦，使得這本書真的誕生。我在 jove 上較在行，幸好他們在唐老鴨的遊樂場較為熟悉。

Andrew S. Tanenbaum

MINIX軟體的適用範圍

從 Prentice-Hall 上得知，MINIX軟體可適用於下列幾種形式：

1. 一群有 640K 隨機存取記憶體 (RAM) 的 IBM PC (或者完全相容者) 磁碟片。這些磁碟片中，含有啓動 (bootable) 系統 (可以被插進 PC 中並執行)，和這系統的原始程式。原始程式磁碟片包含作業系統全部的原始程式，和所有公用程式的原始版本，但 C 語言編譯器除外。(編譯程式的原始程式可由其它方式得到)。另外，這個套裝軟體同樣也可以用在 512K 上，但用戶必須修正一些程式的大小。 ISBN 0-13-583873-8 。
2. 一群為 256K IBM PC 的磁碟片。這個套裝軟體和上述的一樣，只除了 C 語言編譯程式並不包含在內，這是由於記憶體的限制和隨機存取記憶體磁碟機的空間的關係。其它的每個部分，都和 640K 版的一樣。
ISBN 0-13-583881-9
3. 一群 IBM PC / AT (需要 512K) 的磁碟片。這個套裝軟體和 640K PC 的套裝軟體一樣。差別只在隨機存取記憶體磁碟機 (RAM disk) 的空間較為小些，以及用 1.2 M 的磁碟片代替 360K 的磁碟片而已。
ISBN 0-13-583865-7
4. 一個九磁軌、工業用標準的 1600 bpi 以 UNIX “ tar ” 來格式化的磁帶。這個版本包含所有的原始程式，IBM PC 模擬程式，一些程式庫 (Library)，和在 VAX 或其他執行 UNIX 的迷你電腦上執行 MINIX 檔案系統所要用到的一切程式。其實這種移植到其它系統的工作，完全由用戶決定。但由於檔案系統和測試程式都是用 C 寫的，這個步驟應該不難。 ISBN 0-13-583899-1 。

Minix 軟體在上面所述的各種形式中，沒有包含進來的，僅是部分的編譯程式的原始程式。這些程式是用 Amsterdam Compiler Kit (參看 Communications of the ACM, Sept.1983 , pp.654-660) 產生的，這方法也曾用於製造許多其它的語言和機器的編譯程式。它的原始程式可以向下列的公司訂購。

In North and South America:

UniPress Software

2025 Lincoln Highway

Edison, NJ 08817

USA

Telephone : (201)985-8000

In Europe and else where :

Transmediair Utrecht BV

Melkweg 3

3721 RG Bilthoven

Holland

Telephone : (30)78 18 20

目 錄

原序

MINIX軟體的適用範圍

第一章 概論	1 ~ 51
1.1 什麼是作業系統	3
1.1.1 作業系統是個擴充機器	3
1.1.2 作業系統是資源管理者	4
1.2 作業系統的歷史	5
1.2.1 第一代 (1945 ~ 1955)：真空管和插線板	5
1.2.2 第二代 (1955 ~ 1965)：電晶體和批次作業系統	6
1.2.3 第三代 (1965 ~ 1980)：積體電路和多元程式設計	8
1.2.4 第四代 (1980 ~ 1990)：個人電腦	12
1.2.5 MINIX的歷史	13
1.3 作業系統的概念	15
1.3.1 行程	15
1.3.2 檔案	17
1.3.3 shell	21
1.4 系統呼叫	22
1.4.1 行程管理的系統呼叫	26
1.4.2 處理信號的系統呼叫	29
1.4.3 檔案管理的系統呼叫	31
1.4.4 管理目錄的系統呼叫	37
1.4.5 處理保護的系統呼叫	40
1.5 作業系統的結構	42
1.5.1 單一性系統	42
1.5.2 階層式系統	44
1.5.3 虛擬機器	45

1.5.4 顧客 - 侍者模式	47
1.6 本書其餘部份的大綱	49
1.7 摘要	49
問題	50
第二章 行程管理	52 ~ 125
2.1 簡介行程	52
2.1.1 行程模式	52
2.1.2 行程的實做技巧	57
2.2 行程之間的傳輸通訊	59
2.2.1 競爭狀況	59
2.2.2 臨界區域	61
2.2.3 忙碌等待的互斥性	61
2.2.4 睡與醒	66
2.2.5 訊號機	68
2.2.6 事件計數器	71
2.2.7 監督程式	72
2.2.8 訊息傳遞	76
2.2.9 基本運算元的等位關係	82
2.3 典型的行程傳輸問題	86
2.3.1 哲學家晚宴問題	86
2.3.2 讀出者與寫入者問題	90
2.4 行程排班問題	91
2.4.1 循環排班法	93
2.4.2 優先權排班法	94
2.4.3 多重佇列	95
2.4.4 最短的作業最先法	96
2.4.5 政策推導排班法	98
2.4.6 兩段排班法	98
2.5 MINIX 上的行程處理	99

2.5.1 MINIX 的內部結構	100
2.5.2 MINIX 的行程管理	101
2.5.3 MINIX 的行程傳輸處理	102
2.5.4 MINIX 的行程排班問題	103
2.6 MINIX 的行程實做問題	103
2.6.1 MININ 原始碼的結構	104
2.6.2 一般首檔之說明	106
2.6.3 行程之資料結構和首檔說明	109
2.6.4 系統之初設	112
2.6.5 MINIX 的插斷處理	113
2.6.6 核心的公用組合程式碼	116
2.6.7 MINIX 的行程通訊	117
2.6.8 MINIX 的排班技巧	119
2.7 摘要	121
問題	122

第三章 輸入／輸出 126 ~ 215

3.1 輸入／輸出硬體的處理	126
3.1.1 輸入／輸出裝置	127
3.1.2 裝置控制器	128
3.2 輸入／輸出軟體原理	132
3.2.1 輸入／輸出軟體的目標	133
3.2.2 插斷處理程式	134
3.2.3 裝置驅動程式	134
3.2.4 與裝置無關的 I／O 軟體	135
3.2.5 使用者空間的 I／O 軟體	137
3.3 死結	139
3.3.1 資源	140
3.3.2 死結的模式化	140
3.3.3 駝鳥演繹法則	144

3.3.4	發覺與修復	145
3.3.5	預防死結	145
3.3.6	死結之避免	148
3.4	MINIX 的 I/O 概觀	153
3.4.1	MINIX 的插斷控制常式	153
3.4.2	MINIX 的裝置驅動程式	153
3.4.3	MINIX 的與 I/O 裝置無關軟體	157
3.4.4	MINIX 使用者層面的 I/O 軟體	158
3.4.5	MINIX 的死結處理	158
3.5	虛擬磁碟	158
3.5.1	虛擬磁碟的硬體與軟體	159
3.5.2	MINIX 的虛擬磁碟驅動程式概觀	160
3.5.3	MINIX 的虛擬磁碟驅動程式製作	162
3.6	磁碟	163
3.6.1	磁碟的硬體	163
3.6.2	磁碟的軟體	163
3.6.3	MINIX 軟式磁碟驅動程式概觀	169
3.6.4	MINIX 軟式磁碟驅動程式製作	171
3.7	時鐘	174
3.7.1	時鐘硬體	175
3.7.2	時鐘軟體	176
3.7.3	MINIX 時鐘驅動程式概觀	179
3.7.4	MINIX 時鐘驅動程式製作	180
3.8	終端機	182
3.8.1	終端機硬體	182
3.8.2	終端機軟體	186
3.8.3	MINIX 終端機驅動程式概觀	193
3.8.4	MINIX 終端機驅動程式製作	201
3.9	MINIX 的系統作業	206
3.10	摘要	

問題	211
----------	-----

第四章 記憶體管理 216 ~ 279

4.1 最簡記憶體管理法	216
4.1.1 不分頁、置換的單程式作業	216
4.1.2 複程式作業及記憶體	218
4.1.3 固定分割區域的複程式作業	221
4.2 置換	223
4.2.1 變動分割的複程式作業	224
4.2.2 位元圖	226
4.2.3 聯結串列	227
4.2.4 配偶系統	229
4.2.5 置換空間的配置	231
4.2.6 置換系統的分析	231
4.3 虛擬記憶體	232
4.3.1 分頁	233
4.3.2 分段	236
4.4 移頁方法	239
4.4.1 最佳移頁法	239
4.4.2 最近不同移頁法	239
4.4.3 先進先出移頁法	241
4.4.4 最近罕用移頁法	242
4.4.5 LRU 的軟體模擬	243
4.5 分頁系統的設計關鍵	245
4.5.1 工作集合模式	245
4.5.2 局部與通用的配置政策	246
4.5.3 頁的大小	248
4.5.4 實做問題	249
4.6 MINIX 的記憶體管理	251
4.6.1 記憶體的設計	253

4.6.2 訊息處理	255
4.6.3 記憶體管理程式的資料結構和運算法則	256
4.6.4 FORK, EXIT, WAIT 呼叫	259
4.6.5 EXEC 呼叫	260
4.6.6 BRK 呼叫	264
4.6.7 信號處理	264
4.6.8 其它的系統呼叫	266
4.7 MINIX 中記憶體管理實行時的問題	266
4.7.1 首檔	266
4.7.2 主程式	267
4.7.3 FORK, EXIT, WAIT 實行時的問題	268
4.7.4 EXEC 實行時的問題	270
4.7.5 BRK 實行時的問題	271
4.7.6 信號處理實行時的問題	272
4.7.7 其餘系統呼叫實行時的問題	274
4.7.8 記憶體管理程式的應用	274
4.8 摘要	275
問題	276
第五章 檔案系統	280 ~ 379
5.1 使用者眼中的檔案系統	280
5.1.1 檔案基礎	280
5.1.2 目錄	284
5.2 檔案系統設計	286
5.2.1 磁碟空間管理	286
5.2.2 檔案儲存	289
5.2.3 目錄結構	292
5.2.4 共用檔案 (Shared Files)	295
5.2.5 檔案系統可靠性	298
5.2.6 檔案系統的績效	303

5.3 檔案服務程式	306
5.3.1 介面階層 (Interface level)	307
5.3.2 基本更新	308
5.3.3 並行控制 (Concurrency Control)	310
5.3.4 異動作業 (Transaction)	311
5.3.5 複製檔案	313
5.4 安全性	314
5.4.1 安全環境	314
5.4.2 有名的安全缺陷	315
5.4.3 一般的安全侵襲	318
5.4.4 安全性的設計原則	319
5.4.5 使用者識別	319
5.5 保護策略	323
5.5.1 保護定義域	324
5.5.2 取用控制串列	327
5.5.3 資格 (Capability)	329
5.5.4 保護模式 (Protection Models)	331
5.5.5 隱密通道 (Covert channels)	333
5.6 MINIX 檔案系統概觀	335
5.6.1 訊息	335
5.6.2 檔案系統架構	337
5.6.3 位元圖	339
5.6.4 I-nodes	341
5.6.5 快取區塊	342
5.6.6 目錄及路徑	344
5.6.7 檔案描述指標	346
5.6.8 導引管道 (pipe) 及特殊檔案	347
5.6.9 實例：READ 系統呼叫	349
5.7 MINIX 檔案系統的實現	350
5.7.1 標頭檔案 (Header Files)	350