

VAX/VMS操作系统

进程、调度、管理、控制及通讯

国防科技大学研究所六〇二教研室

一九八五年六月

前 言

我们系统地研究和剖析了VAX/VMS操作系统。在此基础上，特陆续地创出一批资料，介绍VMS操作系统。

在这部分资料中，主要介绍VMS操作系统中，有关进程管理，进程调度，进程控制及进程通讯等方面的内容。主要从操作系统内部出发，阐述基本概念，数据结构，及实现的方法和策略。因为这一部分是VMS操作系统最核心的内容之一，所以对许多问题，我们作了较详细的介绍和分析，核心的程序模块也给出了实现的框图或流程图。

文章中错误的地方，请批评指正。

第一章 概 述

VAS/VMS 操作系统是 VAX-11/780 机的核心资源。它是一些控制程序和例行程序的集合。VMS 提供了范围很广的系统服务。系统服务只是一些过程，作为操作系统的一部分，以控制系统资源的合理使用。提供进程间的通讯和控制，实现操作系统的基本功能。

§ 1.1 系统的分层设计

VAX/VMS 操作系统的结构是按层次设计法设计的，它的层次结构见图 1-1。进程的调度和控制在层次结构的核心层。VAX/VMS 系统中的调度程序是使各进程得以合理运行的根本。调度主要负责各进程合理的占用 CPU，控制进程在各状态间的转换等等。

VMS 的核心层由三个相互独立的子系统组成。在各子系统内部是以传统的模块化方法设计的。三个子系统中，虚贮管理子系统和 I/O 子系统均有相应的介绍材料。本部分仅介绍：进程控制子系统的主要内容。

大部分的系统服务（也称广义指令）运行在 VMS 的核方式（K）。VMS 中的系统服务是做为用户进程与操作系统间的界面。外部程序只能通过调用系统服务，才可能得到操作系统的支持。系统服务在 VMS 中的分派控制在中断系统中介绍。详细的说明见相应的介绍材料。记录管理服务（RMS）在执行方式下运行（S）。文件系统的主要内容见“RMS 服务介绍”相应的分派控制也在中断系统中介绍。操作系统的最外层是命令解释程序，命令语言（主

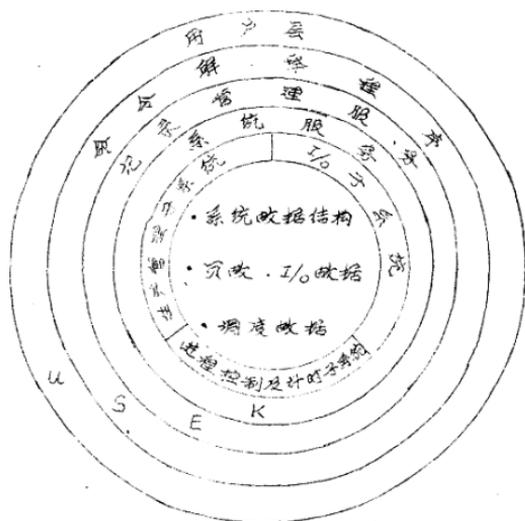


图 1-1 VAX/VMS 操作系统的分层设计

要是 DCL 语言) 是终端用户的系统间的界面。VMS 的命令语言是功能很强的语言, 有关的细节见相应的资料。在用户层主要运行看用户进程和大部分的系统实用程序, 以及库程序。

在介绍这一部分之前, 我们详细阅读了大量的资料, 并用“反汇编器”, 反出并阅读了若干核心的程序模块, 基本搞清了 VMS 的主要控制结构。文章中的框图给的稍粗了些, 主要是为了阅读方便。文章中肯定有不少不成熟的东西, 甚至有许多错误, 望批评指正。

§ 1.2 VMS 提供的访问方式

从图 1-1 可以看到，VAX/VMS 共分了四层。VAX/VMS 分别给各层赋予不同的访问特权。离核心越近，所具有的特权越高。四种访问方式为：

内核方式 (K)：供操作系统核心使用 (包括页面管理，调度和 I/O 子系统)，以及大部分系统服务。

执行方式 (E)：部分系统服务调用和记录管理系统。

管理方式 (S)：用于命令解释之类的服务。

用户方式 (U)：用于用户级程序代码，实用程序，解释程序，调试程序等。

每种访问方式在每进程的进程控制区内均有它自己的堆栈，因此，每个进程可有四个堆栈，每种访问方式一个堆栈，而所以四种访问方式可存在于同一个虚拟地址空间内。

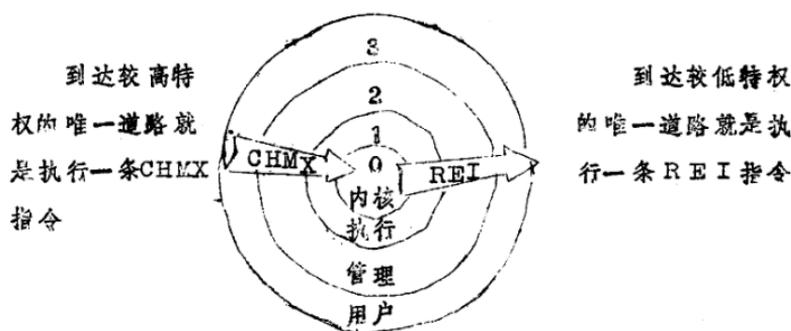


图 1-2 改变访问方式的方法

由于不同的访问方式决定了不同的特权，所以，外层为了得到内层的服务必须改变其访问方式，为此，VAX/VMS规定了改变访问方式的方法，即利用CHMX和REI指令实现各种方式间的转换。

§ 1.3 进程及其有关概念

(1) 进程及进程关联：

• 进程：

所谓进程就是一个具有独立功能的程序关于某个数据集的一次运行活动。在VAX/VMS系统中，一个进程是由系统软件调度的可执行的基本实体，它提供了映像执行的内容，所以进程也是完成用户作业工作的基本实体。当用户向系统注册时，系统就为程序映像的执行建立一个进程。

• 进程关联

一个进程是由虚地址空间及硬件和软件的关联所组成。见图

1-2。

进程的软件关联是由系统对进程进行调度和控制所需的所有数据组成，它主要包括进程的软件优先级、当前的调度状态、进程特权、限额和限制及其它一些信息，如进程名、进程ID等等。

软件关联包括：进程控制块PCB，作业信息块JIB以及进程标题PHD。见附录。

硬件关联的信息主要保存在硬PCB中。而且硬PCB结构本身又是嵌套在进程标题结构中。硬PCB中的域见附录。

在过程关联中，软PCB空间是在建立进程时，动态地从系统

空间中申请，并常驻内存。进程标题位于系统空间的平衡槽中，可被换出内存。

② 关联转换：

当前运行的进程退出运行时，保留进程现场的工作和可执行进程调度执行时恢复现场的工作，称之为关联转换。关联转换主要是涉及硬PCB中的各种寄存器内容。由SVPCTX、LDPCTX指令完成。

③ 映象：

在进程关联中执行的程序称为映象。一个映象由若干过程和数组成。这些数据和过程已由连接程序在一起，即它们已解决了模块间的符号引用，以及分配地址空间。

④ 工作集：

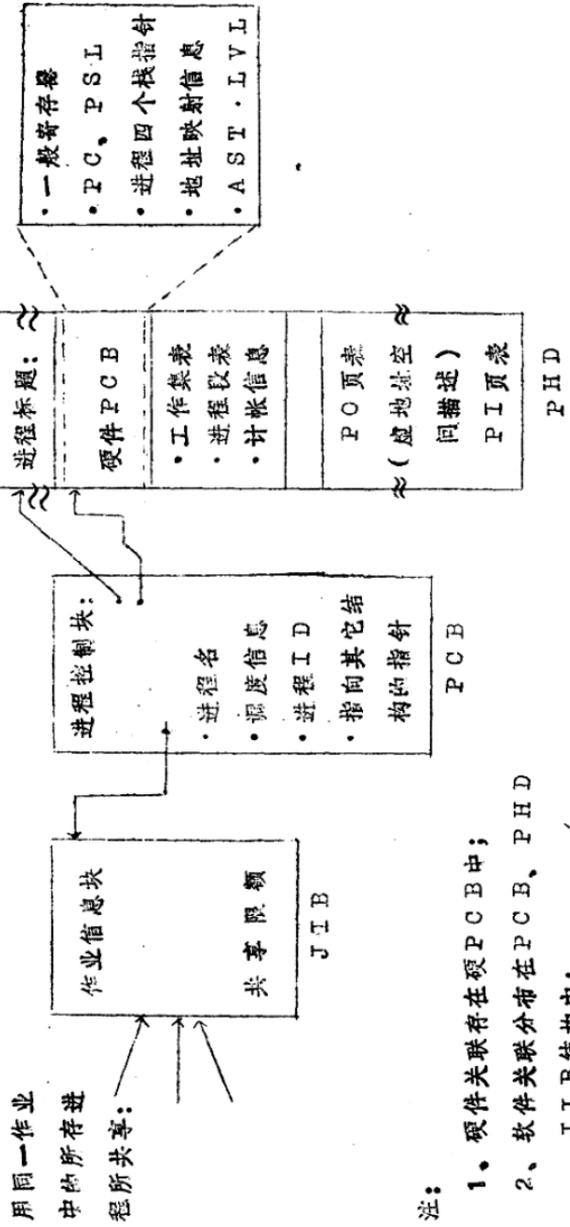
一个进程执行时，该进程虚地址空间页面的一个子集驻留在物理存储器中，这个子集就称为进程的工作集。

⑤ 平衡集：

在任何一个时刻驻留在物理存储器中的所有进程工作集的集合称为平衡集。

⑥ 作业：

作业是完成一项特殊任务的进程族。在进程族中，具有父子关系的一组进程，协同完成一个用户提交的作业。



注:

- 1、硬件关联存在硬PCB中;
- 2、软件关联分布在PCB、PHD JIB结构中;
- 3、“虚地址空间的描述”存储在PO和PI页表中;
- 4、用户程序及控制数据分别占据进程地址空间中的程序区和控制区。

图 1-2 进程关联

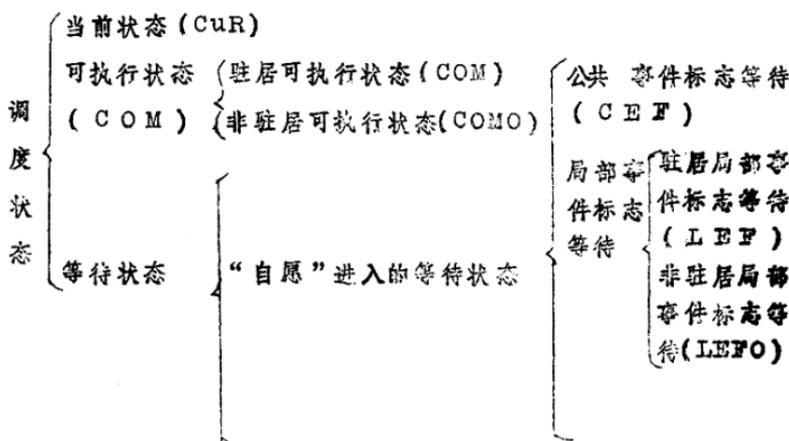
第二章 进程调度

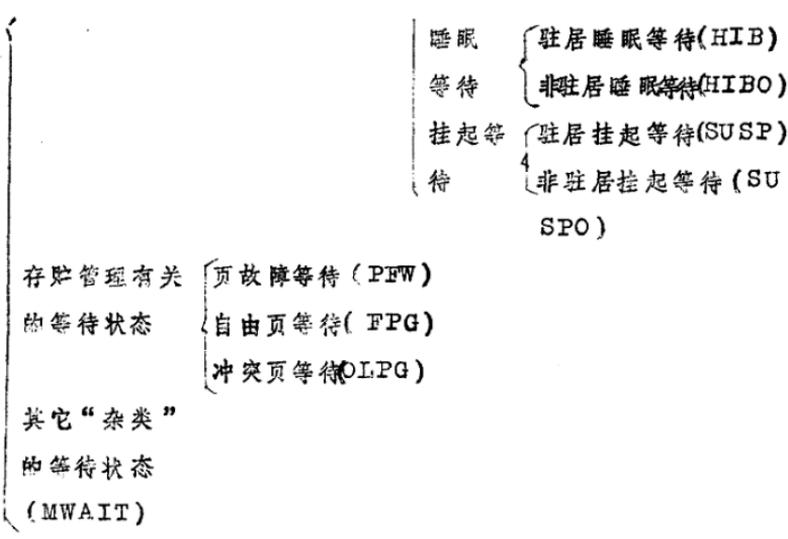
状态及调度控制

所谓进程调度，指的是处理机调度，即“低调”。它主要完成各进程在调度状态间的转换，以及按照选定的策略和算法动态地变化进程的优先权并选择优先权最高的进程在CPU上运行。在多道进程序的系统中调度策略的选择是操作系统的核心问题。VAX/VMS操作系统采用了一种较好的调度策略，充分体现了“分时”和“实时”的特征。本章主要讨论：进程的调度状态；驱动状态变化的系统事件；软件优先权；以及调度程序。

§ 2.1 调度状态及数据结构表示

VMS操作系统将所有进入系统的进程分为如下的调度状态：





以下简单说明各种状态及数据结构表示:

下面所说的数据结构大都是定义在系统虚地址空间中的SDAT模块中。

1. 当前状态 (CUR Current state)

表示在某一特定时刻,正在CPU上运行的进程状态。因为单处理机环境下,处于COR状态的进程最多只能有一个。故数据结构也较简单。SDAT模块中设有一指针: SCH \$ GL-COR PCB指向当前正在运行进程的软PCB地址。

2. 可执行状态 (COM Computable state)

即所谓“就绪状态”,处于该状态的进程已具备了运行条件,只等时机成熟就获得运行。VMS是一个虚拟存贮的操作系统,所以把可执行状态又分为:存贮驻居的可执行状态(Computable

resident state) 记为 COM, 和非存贮驻居的可执行状态 (computable outswapped state) 记为 COMO。

VMS 的调度算法是根据进程优先权的高低选择进程执行的, 而系统中唯有处于这种状态的进程可能被调度运行。为了减少调度中选择最高优先权运行而带来的系统开销, 对 COM 和 COMO 状态的进程设计了精巧的数据结构。

在 SDAT 模块中, 根据 VMS 中的 0~31 个优先级划分, 分别为 COM 和 COMO 状态建立了三十二个队列, 和二个队列概览长字 (QUEUE summary longword)。见图 2-1。

图中表示的是 32 个队列头的结构, 其中每一队列头由四倍字 (quadword) 组成, 表示对应优先级队列的头、尾指针。处于该队列的每一 PCB 用其自身的连接字, 维持这条链。

“队列概览长字”由 32 位组成, 每一位对应一个队列。当某一位为 0 时, 表示对应的队列为空。

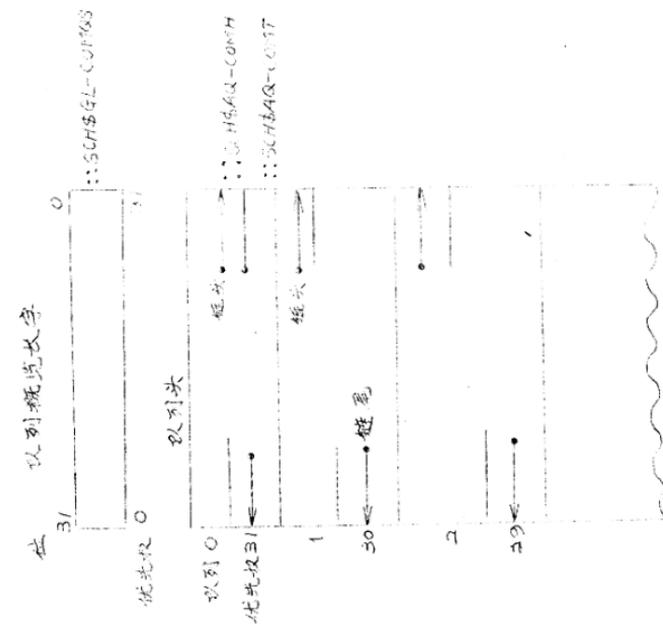
COMO 状态的进程要由交换程序将其换入内存后, 才有可能被调度运行。COM 队列中的进程位是唯一能被调度执行的进程状态。

3、公共事件标志等待 (CEF: Common event Flag wait state)

等待一个或多个公共事件标志的进程到“公共事件块 CEB”上去排队, 公共事件块 CEB 是从系统空间的“非页动态存贮区”中分配的, 队列的头指针 SCH\$GQ-CEBHD 定义在 SDAT 模块中, 有关这一状态的详细实施及数据结构在第四章中讨论。

以下各种等待状态, 它们各自的状态队列具有相同的结构, 每一队列头均由三个四倍字组成。见图 2-2:

COM 状态队列:



COMO 状态队列:

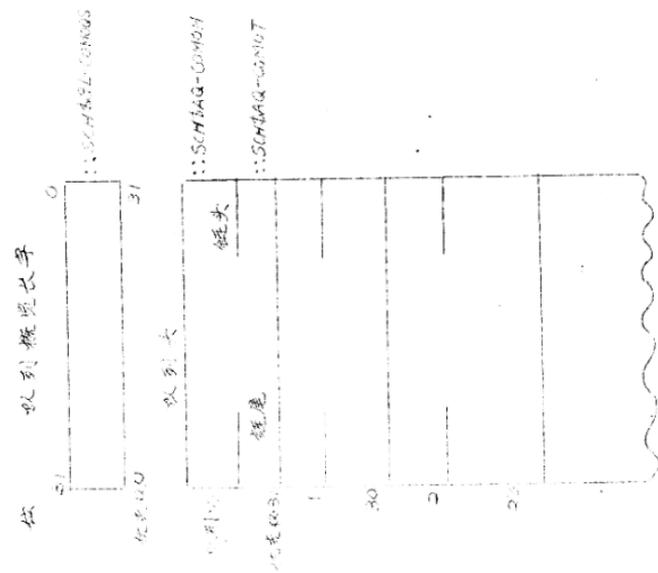


图 2-1 COM COMO 状态队列

等待队列头:

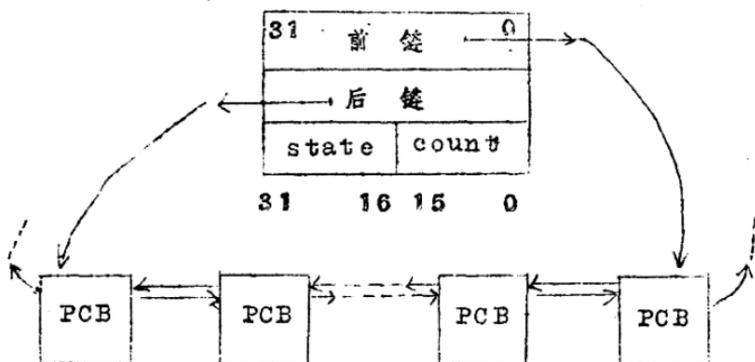


图 2-2 除 CEF 状态外的所有等待状态的队列形成

4. 志愿 (Voluntary) 等待状态:

i) 局部事件等待 LEF: 根据是否驻居于内存分为 LEF LEFO 两个状态。用两条队列表示, 形式如图 2-2。从 CUR 状态变为 LEF 状态。一般是由系统服务: \$WAITER、\$WF LOR \$WFLAND 等引起的;

ii) 睡眠状态 (hibernate): 也分为驻居和非驻居两种: HIB、HIBO。状态队列如图 2-2 “睡眠”系统服务使进程从 CUR 状态转换成 HIB 状态。“\$WAKE 系统服务、ASE 排队。删除进程“夺件”。则能使 HIB 或 HIBO 转换成 Com 或 come 状态。

- iii) 挂起状态 (suspend): 也分为驻居和非驻居两种: SUSP、SUSPO。状态队列如图 2-2。同样, “挂起”系统服务能使进程从 CUR 状态变为 CURO 状态。而“\$RESUME 系统服务和进程删除”, 则能使被挂起的进程变为 com 或 Como 状态。

5. 存贮管理的等待状态:

与存贮管理有关的三个等待状态, 不论是否驻居内存, 均合并在各自己的状态队列中“某一进程被换出(或换入)”时, 只是将该进程软 PCB 中 PCB\$L-STS 域中的 PCB\$V-RES 位清为 0 (或置为 1)。不需移动队列。

- i) 页故障等待状态 (PFM: page fault Wait state) 进程所需的页面不在内存时, 出现页故障。当“故障页面”需要从虚空间读入时, 该进程被放于 PFW 状态。
- ii) 自由页等待状态 (FPG: Free Page Wait state): 当进程需要在自己的工作集中增加一页, 但自由页表已没有可分配的自由页时, 进程进入 FPG 状态。这实际上是一个资源等待。要等待通过“修改的页面写回”, “进程出交换”或“虚拟空间删除”等方法提供了自由页面后, 才可能变为 com 或 como 状态。
- iii) 冲突页等待状态 (COLPG collided page wait state):

出现这种等待一般是因为若干个进程同时出现了“共享页面”的页故障。最初故障的进程进入了EFW状态，而第二个或以后的进程则进入COLPG状态。另外，当某个进程需要一个局部（private）页面，而该页面又已经正在向内存读入时，该进程也进入COLP状态。

当故障页面读入后COLPG状态才可变为com或como状态。

6. 其它杂类的等待状态 (mwait miscellaneous wait state):

MWAIT状态表示进程等待那些未包括在其它等待状态中的资源。这些所有的杂类等待，都链在一条队列中，如图2-2。这些等待的原因，也可分为两类：互斥信号（mutex）等待和资源等待。详细内容列在表2-3中（请见下页）。

i) 资源等待:

资源等待状态是等待那些当时已经用完了的或锁住的某些资源。如非页式的动态存储区已经用完，或公共信箱已经没有空余。等到那些资源重新可用时，处于MWAIT状态的进程才变成com或como状态。

对应资源等待的软PCB中的PCB\$L-EFWm域存在的是：所等待的资源号（\$RSNDEF宏定义的小整数）。

ii) 互斥信号量等待:

等待的原因	PCBSL-EFWM 或内容 (软PCB中)
交互信号等待	数值 (16进制)
系统逻辑	LOG\$AL-mutex 8000 24F4
组名	IOCS\$GL-mutex 8000 24F8
逻辑组	EXES\$GJ-CMBMTX 8000 2620
I/O数据	EXES\$GJ-PGDYNMTX 8000 2624
公共事务	EXES\$GJ-GSDMTX 8000 2628
式动载	EXES\$GJ-SHMGSMTX 8000 2630
局存器	EXES\$GJ-SHMmBMTX 8000 2634
全存器	EXES\$GJ-BNQMNTX 8000 2638
共享存/出	EXES\$GJ-KFIMTX 8000 263C
入队/打印机	UCBS\$L-LP-Mutex 注 1
已知打印机	
资源等待	数值 (16进制)
AST等待 (等待系统或特殊核AST)	RSNS\$-ASTWAIT 00000001
信箱满	RSNS\$-MAILBOX 00000002
非页式动态存储区	RSNS\$-NPDYNMEM 00000003
页式动态存储区	RSNS\$-PGFILE 00000004
页式动态存储 (broadcast) 信息	RSNS\$-PGDYNMEM 00000005
对象激活锁	RSNS\$-BRKTHRU 00000006
作业共享数据 (pooledquota) 等待 (未使用)	RSNS\$-IACLOCK 00000007
	RSNS\$-JQuota 00000008

注：行打印机部件的互斥信号量无固定的地址，信号量地址依赖于该部件所分配的uCB。

互斥信号量是用于保护某些数据结构访问的，若某个进程未获得所申请的信号量使用权，则系统例程便将其置为信号量等待状态，变为 **MUWAIT**。

对应互斥信号量等待状态的软PCB中的PCB\$L-EFLDM域存放所需信号量的系统虚拟地址。

各进程调度状态间的转换关系见图 2-4。

§ 2.2 与调度和调度状态有关的进程关联

3.1 软件进程控制块(软PCB):

SQFL (前链长字)	
SQBL (后链长字)	
PRI (当前优先数)	
PHYPCB 硬PCB物理地址	
STS 软件状态标志	
PRIB 基优先数	STATE 调度状态

详见附录。