

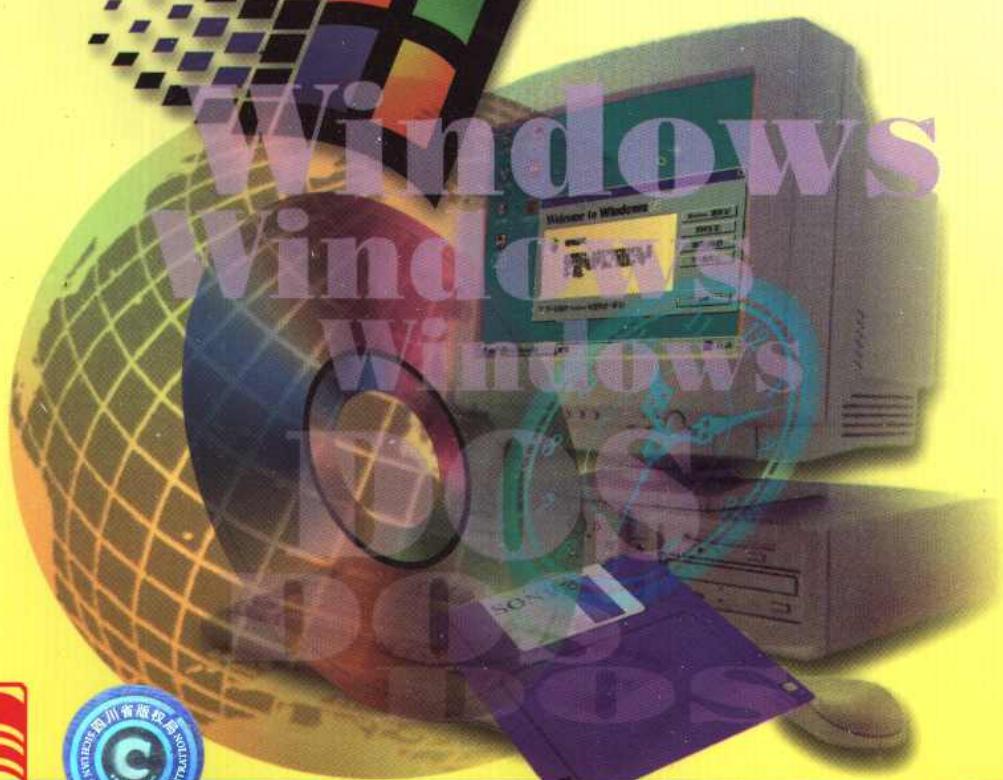
DIANZIKEJIDAXUECHUBANSHE
XILIEJIAOCAI

高等学校
电子信息类 系列教材

大专计算机

操作系统教程

汤子瀛 哲凤屏
汤小丹 王侃雅



电子科技大学出版社

高等学校
电子信息类 系列教材

操作系统教程

汤子瀛 哲凤屏
汤小丹 王侃雅

电子科技大学出版社

前　　言

随着计算机技术的发展，计算机系统的硬件结构日趋复杂，拥有的资源也愈来愈多。为提高资源的利用率、增强系统的处理能力以及使用户能方便和有效地使用计算机，在计算机系统中增设了操作系统（OS）。

操作系统是用以实现上述功能的程序的集合，是配置在计算机硬件上的第一层软件，是对硬件系统的第一次扩充。它在计算机系统中占据了特殊重要的地位，其他所有的软件如汇编程序、编译程序、数据库管理系统等系统软件，以及大量的应用软件，都将依赖于操作系统的支持，取得它的服务。操作系统已成为从大型机直至微机都必须配置的系统软件。

OS 是计算机领域中最活跃的学科之一，其发展极为迅速。为使教材内容能及时地反映时代潮流，我们在本教材中介绍了许多在 90 年代引入或广泛使用的技术，如线程、数据一致性控制技术、系统容错技术、论证技术和数据加密技术等。又因为 20 世纪 90 年代是计算机网络大发展的年代，因而在 90 年代中后期所推出的现代 OS，都毫无例外地提供了面向网络的功能，以方便计算机联网和取得网络所提供的各种服务，故我们在教材中也专门设置了网络操作系统一章。

《操作系统教程》全书共分九章。第一章是 OS 引论，介绍了 OS 的形成、类型、特征和功能等；第二章讲述进程管理，内容包括进程的基本概念、进程的控制、同步、调度和死锁，并介绍了能进一步提高程序并发执行程度的多线程概念；第三章存储管理，主要内容为存储器的连续分配方式和离散分配方式及虚拟存储器等；第四章为设备管理，包括 I/O 控制方式、缓冲管理、设备的分配、设备处理以及磁盘 I/O；第五章是文件管理，主要介绍文件的结构、文件存储空间的管理、目录管理、文件的共享、安全性以及数据一致性的检测；第六章是网络操作系统，其中介绍了网络操作系统的工作模式和在网络环境下 OS 所应具有的功能和服务；第七章是操作系统接口，其中包括当前广泛使用的两种接口；第八章是安全管理，主要介绍为保证计算机系统的安全性，特别是网络环境下系统的安全性所最常使用的各种安全技术；第九章是 UNIX 系统实例，介绍了当前广泛使用的 UNIX 系统的内部结构。

由于得到电子科技大学出版社的大力鼓励和支持，才出版了本教材，在此，特向电子科技大学出版社表示衷心的感谢。限于编者水平，加上时间仓促，本书中必然存在不少错误和不妥之处，恳请读者批评指正。

编者

1999 年 8 月

第一章 操作系统引论

计算机系统由硬件和软件两部分组成，操作系统 OS (Operating System) 是配置在计算机硬件上的第一层软件，是对硬件系统的首次扩充。它在计算机系统中占据了重要的地位，其他所有的软件如汇编程序、编译程序、数据库管理系统以及大量的应用软件，都将依赖于操作系统的支持，取得它的服务。操作系统已成为现代计算机系统（大、中、小微型机）中都必须配置的软件。

在计算机硬件上配置操作系统的主要目标和要求如下：

1. 方便性 操作系统为计算机用户和计算机硬件之间提供接口。用户在 OS 的帮助下能够更方便、快捷、安全、可靠地操纵计算机硬件和运行自己的程序。
2. 有效性 通过 OS 可对计算机硬件和软件资源进行有效地控制和管理，这极大地提高了资源的利用率和改善系统的性能。
3. 可扩充性 由于计算机硬件和体系结构的迅速发展，以及应用的扩大，它们不断地对 OS 提出新的功能和性能要求，因此，要求 OS 必须能方便地扩充和完善，才能满足上述日益增长的需求。
4. 开放性 随着网络化程度的加深，为使出自不同厂家的计算机能通过网络进行集成并能正确、有效地协同工作，实现应用程序的可移植性等，就要求具有统一的开放环境，相应的，OS 必须具有开放性。

1.1 操作系统的形成

从 50 年代至今，在这短短的 40 多年中，OS 取得了如此快速的发展，其推动力主要源于：人们千方百计地提高计算机系统中各种资源的利用率，想方设法地使用户能方便、快捷地使用计算机；此外，计算机器材的不断更新换代、微机的不断发展（由 8 位、16 位发展为 32 位，又进而发展为 64 位）、计算机体系结构的不断发展（由单处理机到多处理机，进而发展到网络和分布式系统），也都推动着 OS 的发展，并产生新的 OS。

1.1.1 人工操作方式

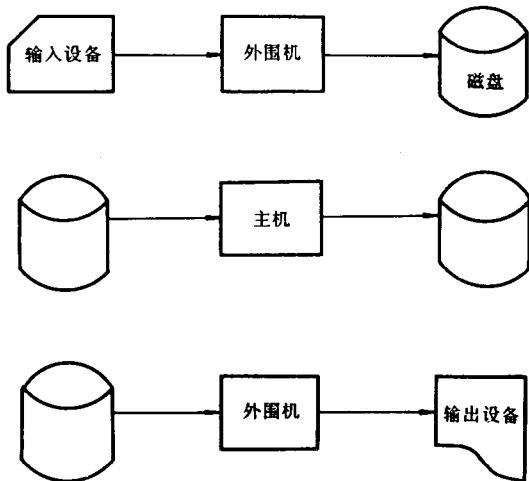
从 1945 年计算机诞生，到 50 年代中期时，还尚未出现 OS，此时仅仅是由用户（即程序员）采用人工操作方式直接使用计算机硬件系统，即由程序员将事先已穿孔（对应于程序和数据）的纸带（卡片），装入纸带输入机（或卡片输入机），再启动该机器将程序和数据输入计算机，然后启动计算机运行。当程序运行完毕并取走计算结果后，才让下一个用户上机。这种人工操作方式有以下两个缺点：

- (1) 用户独占全机 一台计算机的全部资源只能由一个用户独占，致使有些资源空闲。
- (2) CPU 等待人工操作 当用户进行装带（卡）、卸带（卡）等人工操作时，CPU 是空闲的。

可见，人工操作方式产生地降低了计算机资源的利用率，此即所谓的人机矛盾。随着CPU速度的提高、系统规模的扩大，人机矛盾也就变得日趋严重。此外，随着CPU速度的迅速提高而I/O设备的速度却提高得缓慢，又使CPU与I/O设备之间速度不匹配的矛盾更加突出。为了缓和这些矛盾，曾先后出现了通道技术、缓冲技术，但却未能很好地解决上述矛盾，而后来引入的脱机输入输出方式，才获得了较为令人满意的结果。

1.1.2 脱机输入输出(Off-Line I/O)技术

为解决I/O设备的低速问题，在50年代末出现了脱机输入输出技术。该技术是先将用户程序和数据在一台外围机的控制下，预先从低速输入设备输入到磁带或磁盘上。当CPU需要这些数据时，再直接从磁带或磁盘上高速地调入内存。这样就大大地加速了输入过程。类似地，当CPU需要输出时，可立即将输出数据送到磁带或磁盘上，以后再在另一台外围机的控制下，把磁带或磁盘上的处理结果通过相应的输出设备输出，这对程序的执行来说，显然是大大地加速了数据的输出过程。脱机输入输出过程见图1-1所示。由于程序和数据的输入和输出都是在外围机的控制下完成的，或者说它们是脱离主机进行的，故称为脱机输入输出；反之，由主机控制输入输出的方式称为联机输入输出(On-Line I/O)。



脱机输入输出的主要优点如下：

(1) 减少了CPU的空闲时间 装带(卡)、卸带(卡)以及将数据从低速I/O设备送到高速的磁带(盘)上，都是在脱机情况下进行的，它们不占用主机时间，从而有效地减少了CPU的空闲时间，缓和了人机矛盾。

(2) 提高I/O速度 当CPU在运行中需要数据时，是直接从高速的磁带或磁盘上将数据调入内存的，不再是从低速I/O设备上输入，从而大大缓和了CPU和I/O设备速度不匹配的矛盾，进一步减少了CPU的空闲时间。

1.1.3 批处理技术

在早期的脱机I/O方式中，是先把一批作业以脱机输入输出方式输入到磁带上。这意

意味着对作业的处理是成批的，又由于磁带机是顺序存取设备，因而实际上对作业的处理顺序已经排定。为使这一批作业能自动连续地处理，以提高计算机系统的利用率，在系统中还需配置监督程序（Monitor），在它的控制下，先把磁带上的第一个作业装入内存，并将对运行的控制权交给该作业。当该作业处理完后，又把控制权还给监督程序，再由监督程序将第二个作业装入内存……这样自动地一个作业一个作业地进行处理，直到磁带上的所有作业全部完成，这样便形成了早期的批量处理系统，图 1-2 示出了单道批量处理系统的工作流程。

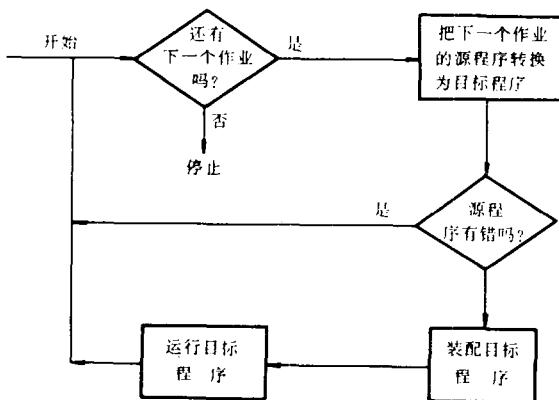


图 1-2 单道批量处理工作流程

由上所述不难看出，批处理技术是在解决人机矛盾和 CPU 与 I/O 设备速率不匹配矛盾的过程中形成的。换言之，批处理技术旨在提高系统资源的利用率和系统吞吐量。但这种单道批处理技术仍然不能很好地利用系统资源，故现已很少使用。

1.1.4 多道程序设计技术

在早期批处理系统中，内存中只存放一道程序，当该程序进行 I/O 操作时，CPU 便处于等待 I/O 完成的空闲状态。为提高 CPU 的利用率和系统的吞吐量，在 60 年代中期又引入了多道程序设计技术，这是同时把几道程序装入内存并允许它们交替执行，共享系统中的各种资源。当正在执行的程序因 I/O 而暂停执行时，CPU 立即转去执行另一道程序；当第二道程序又因 I/O 而暂停执行时，CPU 又转去执行第三道程序。图 1-3 示出了四道程序的运行情况。显然多道程序设计技术提高了 CPU 的利用率，同时也显著地改善了内存和 I/O 设备的利用率，从而也使系统的吞吐量获得大幅度的提高。

多道程序设计技术是一种十分有效又相当复杂的技术，为使在系统中的多道程序之间能协调运行，必须解决下述一系列问题。

1. 处理机管理问题

包括在多道程序之间应如何分配被它们共享的处理机，使 CPU 既能满足各程序运行的需要，又能提高处理机的利用率，以及一旦将处理机分配给某程序后，又应在何时收回等一系列问题。

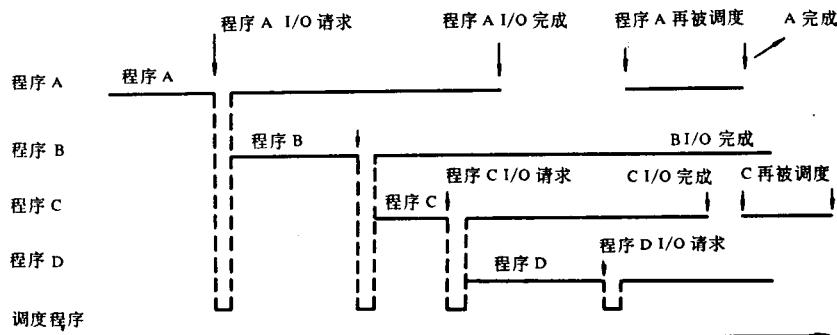


图 1-3 四道程序的运行情况

2. 内存管理问题

包括应如何为每道程序分配必要的内存空间，使它们“各得其所”，且不致因相互重叠而丢失信息，以及应如何防止因某道程序出现异常情况而破坏其他程序等问题。

3. I/O 设备管理问题

系统中可能具有多种类型的 I/O 设备提供给多道程序所共享，应如何分配这些 I/O 设备，以做到既方便用户对设备的使用，又能提高设备的利用率，是 I/O 设备管理应解决的问题。

4. 文件管理问题

在现代的计算机系统中，通常都存放着大量的程序和数据。应如何对它们加以组织以方便用户使用，并保证数据的安全性和一致性，是文件管理中必须解决的问题。

5. 作业管理问题

对于系统中的各种应用程序，其中，有的属于计算型，即以计算为主的程序；有的属于 I/O 型，即以 I/O 为主的程序；又有些作业既重要又紧迫，而有的又要求系统能及时响应，这时，应如何对它们进行组织即作业管理问题。

为此，应在计算机系统中增加一组软件，用以对上述问题进行妥善、有效的处理。这组软件应包括能控制和管理四大资源的软件、合理地对各类作业进行调度的软件以及方便用户使用计算机的软件。正是这样一组软件构成了操作系统。据此，我们可把操作系统定义为：操作系统是一组控制和管理计算机硬件和软件资源、合理地对各类作业进行调度以及方便用户的程序的集合。

1.2 操作系统的基本类型

OS 是配置在计算机上最基本的软件。由于有着形形色色的计算机，相应地，也就有着各种各样的 OS。根据计算机字长的不同，可把 OS 分为 16 位 OS、32 位 OS 和 64 位 OS 等；根据计算机体系结构的不同，可将 OS 分为单处理机 OS、多处理机 OS、网络 OS 等；而根

据 OS 的功能特征，又可将 OS 分为批处理（操作）系统、分时（操作）系统和实时（操作）系统。

1.2.1 批处理系统

基于批处理技术所形成的计算机系统，称为批处理系统。它又可被进一步分成单道批处理系统和多道批处理系统。

一、单道批处理系统 (Simple Batch System)

单道批处理系统是最早出现的一种 OS，严格地说，它只能算作是 OS 的前身而并非是现在人们所理解的 OS。尽管如此，该系统比起人工操作方式已有很大的进步。该系统的主要特征如下：

(1) 单道性 内存中只有一道程序并使之运行，即监督程序每次从磁带上只调入一道程序进入内存运行，仅当该程序完成或发生异常情况时，才调入其后继程序进入内存运行。

(2) 自动性 在顺利的情况下，在磁带上的一批作业会自动地逐个作业依次运行，而无须人工干预。

(3) 顺序性 磁带上的各道作业是顺序地被调入内存，各道作业完成的顺序与它们进入内存的顺序之间，在正常情况下应当完全相同，亦即先调入内存的作业先完成。

二、多道批处理系统

在批处理系统中引入多道程序设计技术后，便形成了多道批处理系统。该系统具有以下特征：

(1) 多道性 在内存中可同时驻留多道程序，并允许它们并发执行，从而有效地提高了资源利用率和系统吞吐量。

(2) 无序性 多个作业完成的先后顺序，与它们进入内存的顺序之间，并无严格的对应关系，即先进入内存的作业可能较后，甚至是最后完成，而后进入内存的作业又可能先完成。

(3) 调度性 作业从被提交给系统开始直至完成，需要经过以下两次调度：

作业调度 这是指按一定的作业调度算法，从外存的后备作业队列中，选择若干个作业调入内存。

进程调度 按一定的进程调度算法，从已在内存的作业中选择一个作业，将处理机分配给它，使之执行。

虽然多道批处理程序诞生于 60 年代，但至今它仍是三大基本操作系统类型之一。在大多数的大、中、小型机中，都配置了它，说明它具有其他 OS 所没有的优点。

多道批处理系统的主要优缺点如下：

1. 主要优点

(1) 资源利用率高 由于在内存中装入了多道程序，使它们共享资源，可保持资源处于忙碌状态，从而使各种资源得到充分利用。

(2) 系统吞吐量大 系统吞吐量是指系统在单位时间内所完成的总工作量。能提高系统吞吐量的原因可归结为：第一，CPU 和其他资源保持“忙碌”状态；第二，仅当作业完

成时，或运行不下去时才进行（作业）切换，系统开销小。

2. 缺点

(1) 平均周转时间长 作业的周转时间是指从作业进入系统开始，直至其完成并退出系统为止所经历的时间。在批处理系统中，由于作业要排队，依次进行处理，因而作业的周转时间较长，通常需几个小时，甚至几天。

(2) 无交互能力 用户一旦把作业提交给系统后直至作业完成，用户都不能与自己的作业进行交互，这对修改和调试程序都是极不方便的。

1.2.2 分时系统 (Time-Sharing System)

分时系统是指在系统的一台主机上，连接了多个带有显示器和键盘的终端，允许多个用户以分时方式共享一台主机的操作系统及计算机系统。即多个用户可通过自己的终端以交互方式使用计算机，共享系统中的硬软件资源。

一、用户需求

如果说，推动多道批处理系统形成和发展的主要动力是提高资源利用率和系统吞吐量，那么，推动分时系统形成和发展的主要动力，则是用户的需求，或者说，分时系统是为了满足用户需求所形成的一种新型 OS。它与多道批处理系统有着截然不同的性能。具体地说，用户的需要表现在以下几个方面。

1. 人—机交互

对于一个程序员来说，每当他编写好一个新的程序时，都需要上机调试。由于新编程序难免有些错误或不当之处需要修改，因而希望能像早期使用计算机一样地“独占全机”，并对它进行直接控制，以便能方便地修改错误，换言之，希望能进行人—机交互。

2. 共享主机

在 60 年代计算机非常昂贵，不可能像现在这样每人独占一台微机，而只能是多个用户共享一台计算机，但用户在用机时应能够像自己独占计算机一样，不仅可以随时与计算机交互，而且应感觉不到其他用户也在使用该计算机。

3. 便于用户上机

对于多道批处理系统，用户上机前必须把自己的作业邮寄或亲自送到机房，这对于用户尤其是远地用户是极不方便的。用户希望能够通过自己的终端直接将作业传送到机器上进行处理，并能对自己的作业进行控制。

二、分时系统实现中的关键问题

实现分时系统最重要的问题是如何使用户能与自己的作业交互，即当用户在自己的终端上键入命令时，系统应能及时接收和及时处理该命令，并将处理结果返回给用户。接着，用户可键入下一条命令，此即人—机交互。应当强调指出，即使有多个用户同时通过自己的键盘键入命令，系统也应能全部地及时接收并及时处理。

要及时接收用户键入的命令或数据并不困难，只需在系统中配置一多路卡。例如，当要在主机上连接 8 个终端时，须配置一个 8 用户的多路卡。多路卡的作用是使主机能“同时”接收用户从各个终端上输入的数据。此外，还必须为每个终端配置一个缓冲区，暂存

用户键入的命令。

人机交互的关键是使用户键入命令后，能及时地控制自己的作业运行或修改自己的作业。为此，各个用户作业都必须在内存中，且应能频繁地获得处理机而运行；否则，用户键入的命令将无法作用到自己的作业上。前面介绍的批处理系统是无法实现人一机交互的，因为通常大多数作业都是驻留在外存上，即使调入内存的作业也经常要经过较长时间的等待后方能运行，因而使用户键入的命令很难及时地作用到自己的作业上。

由此可见，为实现人一机交互，必须彻底地改变原来批处理系统的运行方式。首先，用户作业不能先进入磁盘，然后再调入内存。因为作业在磁盘上不能运行，当然用户也无法与机器交互，因此，作业应直接进入内存。其次，不允许一个作业长期占用处理机，直至它运行结束，或出现 I/O 请求后，方才调度其他作业运行。为此，应该是规定每个作业只运行一很短的时间（例如 0.1 秒钟，通常把这段时间称为时间片），然后便暂停该作业的运行，并立即调度下一个程序运行。如果在不长的时间（如 3 秒）内，能使所有的用户作业都执行一次（一个时间片的时间），便可使每个用户都能及时地与自己的作业交互，从而可使用户的请求得到及时响应。

三、分时系统的实现方法

分时系统的实现有以下三种方法：

1. 单道分时系统

在 60 年代初期，由美国麻省理工学院建立的第一个分时系统 CTSS，便属于单道分时系统。在该系统的内存中只驻留一道程序（作业），其余作业都在外存上。每当内存中的作业运行一个时间片后，便被调至外存（称为调出），再从外存上选一个作业装入内存（称为调入）并运行一个时间片，依此方法使所有的作业都能在一规定的时间内轮流运行一个时间片，这样便使所有的用户都能与自己的作业交互。由于单道分时系统只有一道作业驻留在内存，在多个作业的轮流运行过程中，往往是每个作业都可能被频繁地调进调出多次，开销大，故使系统性能较差。

2. 具有“前台”和“后台”的分时系统

在单道批处理系统中，作业调进调出时，CPU 空闲；内存中的作业在执行 I/O 请求时，CPU 也空闲。为了充分利用 CPU 而引入了“前台”和“后台”的概念（技术）。在具有前、后台的系统中，内存被固定地划分为“前台区”和“后台区”两部分，“前台区”存放按时间片“调进”和“调出”的作业流，“后台区”存放批处理作业。仅当前台调进/调出时，或前台已无作业可运行时，方才运行“后台区”中的作业。

3. 多道分时系统

在分时系统中引入多道程序设计技术后，可在内存中同时存放多道作业，每道程序无固定位置，如果作业都较小，内存中便可多装入几道作业，由系统把已具备运行条件的所有作业排成一个队列，使它们依次轮流地获得一个时间片的处理机来运行。由于切换时作业就在内存，不要花费调入、调出开销，故多道分时系统具有较好的系统性能。现代的分时系统都属于多道分时系统。

四、分时系统的特征

与多道批处理系统比较，分时系统具有下述的特征。

(1) 多路性 允许在一台主机上同时联接多台联机终端，系统按分时原则为每个用户提供服务。宏观上，是多个用户同时工作，共享系统资源；而微观上，则是每个用户作业轮流运行一个时间片。多路性亦即同时性，它提高了资源利用率，从而促进了计算机更广泛的应用。

(2) 独立性 每个用户各占一个终端，彼此独立操作，互不干扰。因此，用户会感觉到是他自己在独占主机。

(3) 及时性 用户的请求能在很短时间内获得响应，此时间间隔是以人们所能接受的等待时间来确定的，通常小于2~3秒钟。

(4) 交互性 用户可通过终端与系统进行广泛的人机对话。其广泛性表现在：用户可以请求系统提供多方面的服务，如文件编辑、数据处理和资源共享等。

1.2.3 实时系统 (Real-Time System)

所谓“实时”，是表示“及时”，而实时系统是指系统能及时（或即时）响应外部事件的请求，在规定的时间内完成对该事件的处理，并控制所有实时任务协调一致的运行。

一、应用需求

虽然多道批处理系统和分时系统，已能获得较为令人满意的资源利用率和响应时间，从而使计算机的应用范围日益扩大，但它们仍然不能满足以下某些应用领域的需要。

1. 实时控制

当把计算机用于生产过程的控制，以形成以计算机为中心的控制系统时，系统要求能实时采集现场数据，并对所采集的数据进行及时处理，进而自动地控制相应的执行机构，使某些（个）参数（如温度、压力、方法等）能按预定的规律变化，以保证产品的质量和提高产量。类似地，也可将计算机用于对武器的控制，如火炮的自动控制系统、飞机的自动驾驶系统，以及导弹的制导系统等。此外，随着大规模集成电路的发展，已制作出各种类型的芯片，并可将这些芯片嵌入到各种仪器和设备中，用来对设备的工作进行实时控制，这就构成了所谓的智能仪器和设备。在这些设备中也需要配置某种类型的、能进行实时控制的系统。通常把用于进行实时控制的系统称为实时控制系统。

2. 实时信息处理

通常，人们把用于对信息进行实时处理的系统称为实时信息处理系统。该系统由一台或多台主机通过通信线路连接到成百上千个远程终端上，计算机接收从远程终端上发来的服务请求，根据用户提出的请求，对信息进行检索和处理，并在很短的时间内为用户作出正确的回答。典型的实时信息处理系统有：飞机或火车票的订票系统、情报检索系统等。

二、实时任务

在实时系统中必然存在着若干个实时任务，这些任务通常与某个（些）外部设备相关，能反映或控制相应的外部设备，因而带有某种程度的紧迫性。可从不同的角度对实时任务

加以分类。

1. 按任务执行时是否呈现周期性来划分

(1) 周期性实时任务 外部设备周期性地发出激励信号给计算机，要求它按指定周期循环执行，以便周期性地控制着某个外部事件。

(2) 非周期性实时任务 外部设备所发出的激励信号，并无明显的周期性，但都必须联系着一个截止时间 (Deadline)。它又可分为：① 开始截止时间：任务在某时间以前必须开始执行；② 完成截止时间：任务在某时间以前必须完成。

2. 根据对截止时间的要求来划分

(1) 硬实时任务 (hard real-time task) 系统必须满足任务对截止时间的要求，否则可能出现难以预测的结果；

(2) 软实时任务 (soft real-time task) 它也联系着一个截止时间，但并不严格，若偶而错过了任务的截止时间，对系统产生的影响也不会太大。

三、实时系统与分时系统特征的比较

我们将从多路性、独立性、及时性、交互性和可靠性五个方面对这两种系统加以比较。

1. 多路性

实时信息处理系统与分时系统一样具有多路性，系统也按分时原则为多个终端用户提供服务；而对实时控制系统，其多路性则主要表现系统在经常对多路的现场信息进行采集，以及对多个对象或多个执行机构进行控制。

2. 独立性

实时信息处理系统与分时系统一样具有独立性。每个终端用户在向实时系统提出服务请求时，是彼此独立地操作，互不干扰；而在实时控制系统中信息的采集和对对象的控制，也都是彼此互不干扰。

3. 及时性

实时信息系统对实时性的要求与分时系统类似，都是以人所能接受的等待时间来确定，而实时控制系统的及时性，则是以控制对象所要求的开始截止时间或完成截止时间来确定的，一般为秒级、百毫秒级直至毫秒级，甚至有的要低于 100 微秒。

4. 交互性

实时信息处理系统虽也具有交互性，但这里人与系统的交互，仅限于访问系统中某些特定的专用服务程序。它不像分时系统那样能向终端用户提供数据处理服务、资源共享等服务。

5. 可靠性

分时系统虽然也要求系统可靠，但相比之下，实时系统则要求系统高度可靠。因为任何差错都可能带来巨大的经济损失，甚至无法预料的灾难性后果。因此，在实时系统中，往往都采取了多级容错措施来保障系统的安全及数据的安全性。

顺便说明，批处理系统、分时系统和实时系统是三种基本的操作系统类型；而一个实际的操作系统，则可能兼有三者或其中两者的功能。例如，在 VAX-11 系列机上所配置的 VMS 操作系统，便是一个兼有分时、实时和批处理功能的操作系统。

1.3 操作系统特征

前面所介绍的三种基本 OS，虽然各有自己的特征，如批处理系统具有成批处理的特征，分时系统具有交互特征、实时系统具有实时特征，但它们也都具有并发、共享、虚拟和异步这四个基本特征。其中，并发特征是 OS 最重要的特征，其他三个特征都是以并发为前提的。

1.3.1 并发 (Concurrency)

并发性是指两个或多个事件在同一时间间隔内发生。在多道程序环境下，并发性是指在一段时间内，宏观上有多道程序在同时运行，但在单处理机系统中，每一时刻却仅能有一道程序执行，故微观上这些程序只能分时地交替执行。倘若在计算机系统中有多个处理器，则这些可以并发执行的程序便可被分配到多个处理器上，实现并行执行，即利用每一个处理器处理一个并发执行的程序。

应当指出，通常的程序是静态实体 (Passive Entity)，它们是不能并发执行的。为使程序能并发执行，系统必须分别为每个程序建立进程。简单说来，进程是指在系统中能独立运行并作为资源分配的基本单位，它是由一组机器指令、数据和堆栈等组成的，是一个活动实体。多个进程之间可以并发执行和交换信息。一个进程在运行时需要一定的资源，如 CPU、存储空间及 I/O 设备等。

在操作系统中引入进程的目的，是使多个程序能并发执行。例如，在一个未引入进程的系统中，在属于一个应用程序的计算程序和 I/O 程序之间，只能是两者顺序执行；但在引入进程后，若分别为计算程序和 I/O 程序各建立一个进程，则这两个进程便可并发执行。由于在系统中具备使计算程序和 I/O 程序同时运行的硬件条件，因而可将系统中的 CPU 和 I/O 设备同时开动起来，实现并行工作，从而有效地提高了系统资源的利用率和系统吞吐量，并改善了系统的性能。

在 OS 中程序的并发执行，将使系统复杂化，以致在系统中必须增设若干新的功能模块，分别用于对处理器、内存、I/O 设备以及文件等系统资源进行管理，并控制系统中作业的运行。事实上，进程和并发是现代 OS 中最重要的基本概念，也是 OS 运行的基础，故我们将在本书第二章中对它们做详细的阐述。

长期以来，进程都是 OS 中可以拥有资源和作为独立运行的基本单位。直到 80 年代中期，人们才又提出了比进程更小的单位——线程。在引入线程的 OS 中，通常仍是把进程作为分配资源的基本单位，而把线程作为独立运行的基本单位。由于线程比进程更小，基本上不拥有系统资源，故而它运行起来更为轻松，能更好地提高系统内多个程序间并发执行的程度。因而，近年来推出的 OS 都引入了线程。

1.3.2 共享 (Sharing)

所谓共享是指系统中的资源可供内存中多个并发执行的进程共同使用。根据资源属性的不同，进程对资源共享的方式也不同，目前主要有以下两种资源共享方式。

1. 互斥共享方式

系统中的某些资源，如打印机、磁带机，虽然它们可以被提供给多个进程使用，但在一段时间内却只允许一个进程访问该资源。当一个进程正在访问该资源时，其他欲访问该资源的进程必须等待。仅当该进程访问完并释放该资源后，才允许另一进程对该资源进行访问。我们把这种资源共享方式称为互斥式共享，而把在一段时间内只允许一个进程访问的资源称为临界资源或独占资源。计算机系统中的大多数物理设备以及某些软件中所用的变量和表格，都属于临界资源，它们要求被互斥地共享。

2. 同时访问方式

系统中还有另一个类资源，允许在一段时间内由多个进程“同时”对它们进行访问。这里所谓的“同时”往往是宏观上的。而在微观上，这些进程可能是交替地对该资源进行访问。典型的可供多个进程“同时”访问的资源是磁盘；一些用重入码编写的文件，也可“同时”共享。

并发和共享是操作系统的两个最基本特征，它们又是互为存在条件。一方面，资源共享是以程序（进程）的并发执行为条件的，若系统不允许程序并发执行，自然不存在资源共享问题；另一方面，若系统不能对资源共享实施有效管理，则也必将影响到程序的并发执行，甚至根本无法并发执行。

1.3.3 虚拟 (Virtual)

操作系统中的所谓“虚拟”，是指通过某种技术把一个物理实体变成若干个逻辑上的对应物。物理实体（前者）是实的，即实际存在的；而后者是虚的，是用户感觉上的东西。相应的，用于实现虚拟的技术，称为虚拟技术。在OS中利用了多种虚拟技术，分别用来实现虚拟处理机、虚拟内存、虚拟外部设备和虚拟信道等。

在虚拟处理机技术中，是通过多道程序设计技术，让多道程序并发执行的方法来分时使用一台处理机，此时，虽然只有一台处理机，但每个终端用户却都认为是有一个CPU在专门为他服务，亦即，利用多道程序技术可以把一台物理上的CPU虚拟为多台逻辑上的CPU，也称为虚拟处理机。

类似地，可以通过虚拟存储器技术，将一台机器的物理存储器变为虚拟存储器，以便从逻辑上来扩充存储器的容量，此时，虽然物理内存的容量可能不大（如4MB），但它可以运行比它大得多的用户程序（如12MB）。这样使用户所感觉到的内存容量要比实际内存容量大得多，认为该机器的内存至少也有12MB。当然，这时用户所感觉到的内存容量是虚的。我们把用户所感觉到的存储器称为虚拟存储器。

我们还可以通过虚拟设备技术，将一台物理I/O设备虚拟为多台逻辑上的I/O设备，并允许每个用户占用一台逻辑上的I/O设备，这样便可使原来仅允许在一段时间内由一个用户访问的设备（即临界资源）变为在一段时间内允许多个用户同时访问的共享设备。例如，原来的打印机属于临界资源，而通过虚拟设备技术可以把它变为多台逻辑上的打印机，供多个用户“同时”打印。此外，也可以把一条物理信道虚拟为多条逻辑信道（虚信道）。在操作系统中虚拟的实现主要是通过分时使用的方法。显然，如果n是某一物理设备所对应的虚拟的逻辑设备数，则虚拟设备的平均速度必然是物理设备速度的1/n。

1.3.4 异步性 (Asynchronism)

在多道程序环境下，允许多个进程并发执行，但只有进程在获得所需的资源后方能执行。在单处理机环境下，系统中只允许一个进程执行，其余进程只能等待。当正在执行的进程提出某种资源请求时，如打印请求，而此时打印机正在为其他某进程打印，这时，正在执行的进程必须等待且放弃处理机，直到打印机空闲，并又把处理机分配给该进程时，该进程方能继续执行。可见，由于资源等因素的限制，通常，进程的执行并非是“一气呵成”，而是以“走走停停”的方式运行。

内存中的每个进程在何时能获得处理机而执行，何时又会因提出某种资源请求而暂停，以及进程以怎样的速度向前推进，每道程序总共需多少时间才能完成，都是不可预知的。由于各用户程序性能的不同，如有的侧重于计算而较少需要 I/O；而又有的程序其计算少而 I/O 多，这样，很可能是先进入内存的作业后完成；而后进入内存的作业先完成。或者说，进程是以人们不可预知的速度向前推进，此即进程的异步性。尽管如此，但只要运行环境相同，作业经多次运行，都会获得完全相同的结果。因此，异步运行方式是允许的，是操作系统的一个重要特征。

1.4 操作系统的任务和功能

操作系统的的主要任务是为多道程序运行提供良好的运行环境，以保证多道程序能有条不紊地、高效地运行并能最大程度地提高系统中各种资源的利用率和方便用户的使用。为实现上述任务，操作系统应具有这样几方面的功能：处理机管理、存储器管理、设备管理和文件管理。为了方便用户使用操作系统，还须向用户提供一个使用方便的用户接口。此外，当今的网络已相当普及，已有愈来愈多的计算机接入网络中，为了方便计算机联网，在 OS 中已增加了面向网络的服务和功能。

1.4.1 处理机管理的任务和功能

在传统的多道程序系统中，处理机的分配和运行都是以进程为基本单位，因而对处理机的管理可归结为对进程的管理。处理机管理的主要功能是创建和撤消进程，对诸进程的运行进行协调、实现进程之间的信息交换，以及按照一定的算法把处理机分配给进程。

一、进程控制

在传统的多道程序环境下，要使作业运行，必须先为它创建一个或几个进程，并为之分配必要的资源。当进程运行结束时，立即撤消该进程，以便及时回收该进程所占用的各类资源。进程控制的主要任务是为作业创建进程、撤消已结束的进程，以及控制进程在运行过程中的状态转换。

二、进程同步

为使多个进程能有条不紊地运行，系统中必须设置进程同步机制。进程同步的主要任务是对多个进程的运行进行协调。有两种协调方式：(1) 进程互斥方式：这是指诸进程在

对临界资源进行访问时，应采用互斥方式；（2）进程同步方式：指在相互合作完成共同任务的诸进程间，由同步机构对它们的执行次序加以协调。

为了实现进程同步，系统中必须设置进程同步机制。最简单的用于实现进程互斥的机制，是为每一个临界资源配置一把锁 W，当锁已打开时，进程可以对该临界资源进行访问；而当锁已关上时，则禁止进程访问该临界资源；而实现进程同步最常用的机制则是信号量，我们将在第二章中做详细的介绍。

三、进程通信

在多道程序环境下，为了加速应用程序的运行，应在系统中建立多个进程，这些进程相互合作去完成一共同任务，而在这些进程之间，又往往需要交换信息。例如，有三个相互合作的进程，它们是输入进程、计算进程和打印进程。输入进程负责将所输入的数据传送给计算进程；计算进程利用输入数据进行计算，并把计算结果传送给打印进程，最后由打印进程把计算结果打印出来。进程通信的任务就是用来实现在相互合作进程之间的信息交换。

当相互合作的进程处于同一计算机系统时，通常在它们之间是采用直接通信方式，即由源进程利用发送命令直接将消息（Message）挂到目标进程的消息队列上，以后由目标进程利用接收命令从其消息队列中取出消息。

四、调度

在后备队列上等待的每个作业，通常都要经过调度才能执行，其中包括作业调度和进程调度两步。作业调度的基本任务，是从后备队列中按照一定的算法，选择出若干个作业，为它们分配其必需的资源（首先是分配内存）。在将它们调入内存后，便分别为它们建立进程，使它们都成为可能获得处理机的就绪进程；并将它们按照一定的算法插入就绪队列。而进程调度的任务，则是从进程的就绪队列中，按照一定的算法选出一新进程，把处理机分配给它，并为它设置运行现场，使进程投入执行。

1.4.2 存储器管理的任务和功能

存储器管理的主要任务，是为每道程序分配内存空间，确保每道程序在规定的空间中运行，方便用户使用存储器，提高存储器的利用率，以及能从逻辑上来扩充内存。为此，存储器管理应具有以下功能：内存分配、内存保护、地址映射和扩充内存等。

一、内存分配

内存分配的主要任务是为每道程序分配内存空间，使它们“各得其所”；提高存储器的利用率，以减少不可用的内存空间；允许正在运行的程序申请附加的内存空间，以适应程序和数据动态增长的需要。

为了实现内存分配，在实现内存分配的机制中应具有以下的数据结构和功能：（1）内存分配的数据结构：该结构用于记录内存空间的使用情况，作为内存分配的依据；（2）内存分配功能：系统按照一定的内存分配算法为用户程序分配内存空间；（3）内存回收功能：系统对用户已不再需要的内存，通过用户的释放请求去完成系统的回收功能。

二、内存保护

内存保护的主要任务，是确保每道用户程序都在自己的内存空间中运行，互不干扰。进一步说，绝不允许用户程序访问操作系统的程序和数据；也不允许转移到非共享的其他用户程序中去执行。为此，计算机必须提供必要的硬件保护机制，并用软件相配合实现内存保护功能。

一种比较简单的内存保护机制，是设置两个界限寄存器，分别用于存放正在执行程序的上界和下界。系统须对每条指令所访问的地址进行越界检查，如果发生越界，便发出越界中断信号，以停止该程序的执行。如果这种检查完全由软件实现，则每执行一条指令，便需要增加若干条指令去进行越界检查，这将显著降低程序的运行速度。因此，越界检查都由硬件实现。当然，对发生越界后的处理，还须与软件配合来完成。

三、地址映射

一个应用程序（源程序）经编译和执行后，所形成的地址范围称为“地址空间”，其中的地址称为“逻辑地址”或“相对地址”。此外，由内存中的一系列单元所限定的地址范围称为“内存空间”，其中的地址称为“物理地址”。在多道程序环境下，地址空间中的逻辑地址和内存空间中的物理地址，是不可能一一对应的。因此，存储器管理必须提供地址映射功能，以将地址空间中的逻辑地址转换为内存空间中与之对应的物理地址。该功能通常是在地址映射表的支持下完成的。

四、内存扩充

尽管配置在计算机中的内存愈来愈大，从数十千字节扩大到数兆字节，甚至更大，但仍然难于满足用户的需要。在存储器管理中的“内存扩充”任务，并非是去增加物理内存的容量，而是借助于虚拟存储技术，从逻辑上去扩充内存容量，使用户所感觉到的内存容量比实际内存容量大得多。换言之，它是使系统能运行的应用程序，其所要求的内存容量可以比物理内存大得多；或者是让更多的用户程序能并发运行。这样，既满足了用户的需要，改善了系统性能，又基本上不增加硬件投资。

1.4.3 设备管理的任务和功能

设备管理的主要任务，是完成进程提出的 I/O 请求，为进程分配 I/O 设备；提高 CPU 和 I/O 设备的利用率、提高 I/O 速度，以及方便用户使用 I/O 设备。为了实现上述任务，设备管理应具有缓冲管理、设备分配和设备处理，以及虚拟设备等功能。

一、缓冲管理

利用缓冲区可以有效地缓和 CPU 和 I/O 设备速度不匹配的矛盾，提高 CPU 和 I/O 设备的利用率。因此，在现代计算机系统中都设置了多种类型的缓冲，由 OS 中的缓冲管理机构把它们管理起来。

缓冲管理的主要功能是管理好各种类型的缓冲区，如字符缓冲区和字符块缓冲区，以缓和 CPU 和 I/O 速度不匹配的矛盾，最终达到提高 CPU 和 I/O 设备利用率，进而提高系