

并行与分布计算技术丛书

并行操作系统原理与技术

The Principle and Technique of Parallel Operating System

夏卫民 罗 宇 吴庆波 著
俞 晓 肖 依 周良源

国防工业出版社

·北京·

图书在版编目(CIP)数据

并行操作系统原理与技术/夏卫民等著. —北京:国防工业出版社, 2002. 2
(并行与分布计算技术丛书)
ISBN 7-118-02670-0

I . 并... II . 夏... III . 并行控制系统: 操作系统(软件) IV . TP316

中国版本图书馆 CIP 数据核字(2001)第 073311 号

国防工业出版社出版发行

(北京市海淀区紫竹院南路 23 号)

(邮政编码 100044)

三河市腾飞胶印厂印刷

新华书店经售

开本 787×1092 1/16 印张 14^{3/4} 320 千字

2002 年 2 月第 1 版 2002 年 2 月北京第 1 次印刷

印数: 1~4000 册 定价: 36.00 元

(本书如有印装错误, 我社负责调换)

致 读 者

本书由国防科技图书出版基金资助出版。

国防科技图书出版工作是国防科技事业的一个重要方面。优秀的国防科技图书既是国防科技成果的一部分,又是国防科技水平的重要标志。为了促进国防科技和武器装备建设事业的发展,加强社会主义物质文明和精神文明建设,培养优秀科技人才,确保国防科技优秀图书的出版,原国防科工委于1988年初决定每年拨出专款,设立国防科技图书出版基金,成立评审委员会,扶持、审定出版国防科技优秀图书。

国防科技图书出版基金资助的对象是:

1. 在国防科学技术领域中,学术水平高,内容有创见,在学科上居领先地位的基础科学理论图书;在工程技术理论方面有突破的应用科学专著。
2. 学术思想新颖,内容具体、实用,对国防科技和武器装备发展具有较大推动作用的专著;密切结合国防现代化和武器装备现代化需要的高新技术内容的专著。
3. 有重要发展前景和有重大开拓使用价值,密切结合国防现代化和武器装备现代化需要的新工艺、新材料内容的专著。
4. 填补目前我国科技领域空白并具有军事应用前景的薄弱学科和边缘学科的科技图书。

国防科技图书出版基金评审委员会在总装备部的领导下开展工作,负责掌握出版基金的使用方向,评审受理的图书选题,决定资助的图书选题和资助金额,以及决定中断或取消资助等。经评审给予资助的图书,由总装备部国防工业出版社列选出版。

国防科技事业已经取得了举世瞩目的成就。国防科技图书承担着记载和弘扬这些成就,积累和传播科技知识的使命。在改革开放的新形势下,原国防科工委率先设立出版基金,扶持出版科技图书,这是一项具有深远意义的创举。此举势必促使国防科技图书的出版随着国防科技事业的发展更加兴旺。

设立出版基金是一件新生事物,是对出版工作的一项改革。因而,评审工作需要不断地摸索、认真地总结和及时地改进,这样,才能使有限的基金发挥出巨大的效能。评审工作更需要国防科技和武器装备建设战线广大科技工作者、专家、教授,以及社会各界朋友的热情支持。

让我们携起手来,为祖国昌盛、科技腾飞、出版繁荣而共同奋斗!

**国防科技图书出版基金
评审委员会**

国防科技图书出版基金 第三届评审委员会组成人员

名誉主任委员 怀国模

主任委员 黄 宁

副主任委员 殷鹤龄 高景德 陈芳允 曾 锋

秘书长 崔士义

委员 于景元 王小谟 尤子平 冯允成

(以姓氏笔划为序)

刘 仁 朱森元 朵英贤 宋家树

杨星豪 吴有生 何庆芝 何国伟

何新贵 张立同 张汝果 张均武

张涵信 陈火旺 范学虹 柯有安

侯正明 莫梧生 崔尔杰

并行与分布计算技术丛书编委会

主 编 卢锡城

副 主 编 周兴铭 汪成为 李国杰

主 审 陈火旺

编委委员 施伯乐 康继昌 尤晋元 康立山

沈隆均 李晓梅 朱传琪 王志英

杨学军 杨晓东 李思昆 王怀民

常务秘书 肖 政

总序

并行与分布计算技术是实现高性能计算的重要技术途径。高性能计算机技术是现代科学研究、工程技术开发和大规模数据处理的关键支撑技术。没有高性能计算机，大量复杂问题的计算和事务处理就无法在合理的时间内完成。利用高性能计算机还可以解决一些仅靠理论及实验方法无法解决的问题，分析处理靠传统技术无法应付的海量数据，例如核爆炸模拟、宇宙的形成及演变过程研究、中长期天气预报、石油地质勘探、数字地球和大规模事务处理等。人类对高性能计算能力的需求永无止境。计算速度、存储容量、通信带宽是衡量高性能计算能力的重要技术指标。一些推动人类文明和社会信息化的重大挑战性问题需要百万亿次每秒、千万亿次每秒以上的计算速度，需要万亿字节以上的存储容量和万亿位每秒以上的通信能力。人们已经认识到，高性能计算与网络通信技术是战略技术，是科技创新体系的基础技术，是反映一个国家综合实力的重要标志之一。

由于种种因素的制约及使用计算机方式的改变，在单台计算机系统的性能与功能难以满足应用需求的情况下，由多个计算节点经专门设计的互联网络紧密耦合而构成的大规模并行处理计算机系统，以及由多个自主的高性能计算节点经计算机网络连接组成的分布处理环境，将成为高性能计算机领域的两大重要研究方向。并行计算和分布计算是既有区别又联系密切的两个概念，前者重在发掘计算过程中的并行性，后者则重在有效组织管理各类异构资源，挖掘功能上的并发性。随着基于先进计算机网络的分布并行计算概念的发展，并行计算与分布计算在很多应用领域正面临共同的追求目标和技术挑战，如高效、实用的计算模型、计算方法、并行机制等。

20世纪90年代以来，随着高性能计算和网络计算技术的普及，并行与分布计算技术正渗透到现代社会的各个领域，用户通过高性能网络使用各类计算资源提供的服务已成为信息化社会中计算机应用的一种重要形式。事实上，各类基于并行与分布计算的应用系统正在工业、交通、金融、科研、政府、国防等部门支撑着现代社会的高效运行。

20多年来，我国科技人员依靠自己的力量，勇于开拓，奋勇拼搏，研制成功了多种型号序列的高性能计算机系统。国防科学技术大学计算机学院是我国研制高性能计算机系统的重要基地，多种巨型计算机的研制成功及推广应用，打破了国外对我国的技术封锁，缩短了我国同发达国家技术水平上的差距，为推动我国高性能计算机及并行与分布计算技术的发展作出了重要贡献。

为促进并行与分布计算技术领域的研究，国防工业出版社组织国防科技大学计算机学院有关专家、教授撰著了本套丛书。本丛书以并行与分布计算机系统组成为纲，结合作者多年科研与工程实践以及当前研究的热点问题，涵盖了计算机系统结构、计算机网络、系统软件、应用软件、计算方法等多方面内容。本丛书由以下9本专著组成：《并行计算机体系结构技术》论述并行计算机研究和设计的理论及工程问题；《先进计算机网络技术》论

述高性能网络计算的各种关键技术;《并行操作系统原理与技术》论述并行操作系统的机制、基本原理和主要实现技术;《并行编译方法》论述并行编译系统理论、编译器的设计方法;《分布计算——网络化软件新技术》论述分布计算技术的基本概念和关键技术;《分布式数据库技术》论述统一逻辑分布式数据库技术;《数字系统并行 CAD 技术》论述数字系统并行 CAD 的理论和技术;《可扩展并行算法的设计与分析》论述可扩展并行算法的设计与分析的理论和方法;《并行与分布式可视化技术及应用》论述并行与分布式可视化的各种技术及其应用。本丛书既介绍了当前国际上该领域的最新技术发展,又汇集了作者多年的研究成果和工程经验。丛书注重可读性,适合从事该领域工作和学习的科技人员、高等院校高年级学生及研究生作为工作和学习的参考书。

本丛书被列入“九五国家重点图书选题规划”,并获得国防科技出版基金的资助。愿本丛书的出版能为并行与分布计算技术研究园地增添一朵花絮,为以后的研究工作提供有价值的参考。因时间和能力所限,书中不足之处,恳请读者指正。

并行与分布计算技术丛书编委会

前　　言

并行处理是当今计算机发展的一个热点,自 20 世纪 90 年代以来,信息时代已进入了并行处理的年代,尽管并行处理技术还有待进一步发展和提高,但它已被广泛地应用于当今的科学计算、信息处理和移动通信,特别是高性能计算等领域。继“银河 - I”亿次向量巨型机于 1983 年诞生之后,我国已先后成功研制出“银河 - II”、“银河 - III”、“神州 - I”、“神州 - II”、“曙光 1000”和“曙光 2000”等多个并行计算机系列。并行操作系统是并行计算机的重要组成部分,它也随着并行处理技术和并行计算机体系结构的发展而发展,并且越来越受到人们的重视。现在我国已有大量的科研人员在从事并行处理方面的研究工作,大专院校的学生们也在进行或将要进行并行处理方面的学习和探索工作,但是目前在并行处理方面的中文书籍和参考资料却很缺乏,特别是并行操作系统方面的书籍还没见到。为此,我们在参阅了大量国内外有关并行处理资料的基础上,从分析传统操作系统与并行处理不相适应的地方入手,并结合我们自己的实践,撰写了本书,希望能起到抛砖引玉之功效。

本书共 8 章,论述了并行操作系统的并行机制、基本原理和主要技术。

第 1 章概论。简要介绍了并行操作系统的发展历程,展示了并行操作系统随现代科技进步和体系结构的发展而进化的艰苦历程,给出了现代并行操作系统的主要结构、特点和适应面。

第 2 章进程与线程。介绍了进程与线程的控制模型、运行机理以及它们的状态演变过程,特别介绍了在并行操作系统中广泛使用的线程的工作原理与实现方法。

第 3 章处理机调度。在分析一般进程调度的基础上,介绍了多机调度的基本概念以及线程调度和实时调度。

第 4 章存储管理。首先介绍了现代计算机存储体系结构,分析了存储体系中的层次关系及分布模式,针对不同的存储器模型介绍了用户空间的组织方法、数据访问模式,特别介绍了分布存储器系统上数据的迁移技术和数据一致性管理技术。

第 5 章同步与互斥。首先分析了同步机制在并行机系统中的重要地位,分析了单机系统中同步机制的不足,介绍了多机系统中常见的同步手段,特别对不同种类的锁机制进行了解剖和分析。

第 6 章机间通信技术。介绍了共享存储器和分布存储器环境下消息的发送与接收,特别介绍了现代流行的端口技术、主动消息及消息传递接口界面。

第 7 章并行 I/O 和并行文件系统。重点讨论了并行 I/O、分布式文件系统、并行文件系统的技术和原理,并给出了相关的实现实例。

第 8 章 Mach 操作系统核心分析。在分析源代码的基础上,对 Mach 微内核的体系结构、进程管理、存储及虚存管理以及进程间通信进行了一定的剖析。

本书第1章由周良源编写,第2章、第3章由罗宇编写,第4章由俞晓和夏卫民编写,第5章由夏卫民和周良源编写,第6章由肖依编写,第7章由吴庆波编写,第8章由刘忠编写。卢宇彤、陈松政和郑立刚等同志参加了本书的撰写工作。全书由夏卫民统稿。杨学军教授在百忙中仔细审阅了全书并提出了许多指导意见。在此,对关心和帮助过本书编著工作的所有同志表示由衷的感谢。

限于编者的水平,书中疏漏及错误之处恳请专家、读者批评指正。

目 录

第 1 章 概述	1
1.1 并行操作系统的发展史	3
1.2 多处理机系统结构	6
1.3 并行操作系统的分类	9
1.3.1 按控制方式	9
1.3.2 内核结构	10
第 2 章 进程与线程	14
2.1 进程描述	14
2.1.1 什么是进程	15
2.1.2 操作系统控制结构.....	15
2.1.3 进程控制结构	17
2.2 进程状态	19
2.2.1 进程的创建与结束.....	20
2.2.2 进程状态变化模型.....	21
2.2.3 进程挂起	23
2.3 进程控制	25
2.3.1 执行模式	25
2.3.2 进程切换	26
2.3.3 操作系统运行模型.....	27
2.4 多处理机与线程	30
2.4.1 对称多处理机	30
2.4.2 进程模型和线程模型	32
2.4.3 线程实现	40
第 3 章 处理机调度	45
3.1 进程调度	45
3.1.1 分级调度	46
3.1.2 进程调度方式与实现	47
3.1.3 调度算法	49
3.2 多机调度设计	53
3.2.1 并行系统结构与调度	54
3.2.2 线程调度	57
3.3 分布主存多机分配和调度	62

3.3.1 分布主存并行计算机的类型	62
3.3.2 处理器分配和调度算法概述	63
3.3.3 典型的分配和调度算法示例	66
3.4 典型并行操作系统处理器调度分析	70
3.4.1 LWP、纯内核线程和用户级线程	70
3.4.2 UNIX SVR4.2MP 的处理机调度	71
3.4.3 Windows NT 处理机调度	73
第 4 章 存储管理	75
4.1 存储器的结构模式	75
4.1.1 单一访问模式	75
4.1.2 非一致存储访问模式	76
4.1.3 非远程存储访问模式	77
4.1.4 全高速缓存存储访问模式	77
4.2 虚拟存储	78
4.2.1 地址空间	78
4.2.2 地址转换	79
4.2.3 UMA 模式的存储管理	81
4.2.4 其它模式的存储管理	81
4.3 数据一致性的管理	82
4.3.1 一致性的模式	82
4.4 存储器一致性模型与编程模型	88
4.4.1 Cache 一致性协议	89
4.5 虚共享系统与共享存储空间编程模型	90
4.6 并行环境下的数据访问与数据局部性	92
4.7 页迁移技术	92
4.7.1 页迁移的 4 个关键问题	93
4.7.2 页故障触发的页迁移技术	95
4.7.3 Cache 失效触发的页迁移技术	97
4.7.4 懒宿主页迁移	100
4.7.5 3 种页迁移技术的比较	103
第 5 章 同步与互斥	105
5.1 同步的提出	105
5.1.1 几个单机系统中常见的典型例子	105
5.2 单机系统中常采取的主要方法	106
5.3 并行处理环境下同步与互斥的特殊性	109
5.3.1 数据保护问题	111
5.3.2 锁语义的修改	111
5.3.3 中断屏蔽问题	112
5.3.4 信号量缺陷	112

5.3.5 丢失唤醒问题	113
5.3.6 极度集中问题	113
5.3.7 共享与互斥访问	113
5.4 共享主存多处理机的同步与互斥	113
5.4.1 共享主存多机系统环境下常见的基本操作	114
5.4.2 操作系统环境下常见的同步与互斥操作	116
5.4.3 多处理机系统中几个值得讨论的问题	117
5.5 分布主存环境下的同步与互斥	119
5.5.1 集中式控制方法	119
5.5.2 无仲裁者竞争方法	119
5.5.3 令牌环方法	120
5.5.4 失效时的解决方法	120
5.6 一些实例	120
5.6.1 UNIX System V Release 4	121
5.6.2 其它实现	122
第6章 机间通信技术	126
6.1 基本概念	126
6.1.1 消息传递	126
6.1.2 共享变量	127
6.2 通信协议	128
6.2.1 分层网络的体系结构	128
6.2.2 ISO/OSI 7 层通信参考模型	132
6.2.3 TCP/IP 网络参考模型	133
6.3 MPP 操作系统通信技术	135
6.3.1 PVM 并行虚拟机系统	135
6.3.2 消息传递标准 MPI	139
6.4 主动消息通信及其应用编程界面	142
6.4.1 简介	142
6.4.2 主动消息 API: 端点和群	145
6.4.3 主动消息 API: 并发和同步	150
6.4.4 主动消息 API: 传输操作	154
6.4.5 Solaris/Myrinet 端点管理	157
6.5 多通信方法技术和 NEXUS 实现	158
6.5.1 原因	158
6.5.2 多方法通信	159
6.5.3 NEXUS 实现方法	159
第7章 并行 I/O 和并行文件系统	164
7.1 分布式文件系统	165
7.1.1 分布式文件系统的设计	166

7.1.2 分布式文件系统的实现	168
7.1.3 分布式文件系统的新技术	172
7.1.4 分布式文件系统举例	173
7.2 并行 I/O 和并行文件系统	181
7.2.1 文件抽象	181
7.2.2 并行 I/O 的方法	183
7.2.3 并行文件系统	184
7.2.4 Galley 并行文件系统	184
7.2.5 PVFS 并行文件系统	191
7.2.6 Panda 并行 I/O 库	194
7.2.7 MPI - IO	196
第 8 章 Mach 操作系统核心分析	200
8.1 Mach 概述	200
8.2 Mach 微内核	201
8.2.1 微内核	201
8.2.2 Mach 微内核中的对象	202
8.2.3 客户/服务器模型	202
8.3 Mach 的进程管理	203
8.3.1 任务	203
8.3.2 线程	203
8.3.3 Mach 的线程调度	205
8.4 Mach 的存储管理	206
8.4.1 存储对象	207
8.4.2 Mach 的虚存对象	207
8.4.3 存储对象到虚空间的映射	208
8.5 Mach 的通信机制	209
8.5.1 基本概念	209
8.5.2 端口权限保护机制的实现	209
8.5.3 实现端口权限保护的一般方法	212
8.5.4 网络通信	213
参考文献	214

Content

Chapter 1 Summary	1
1.1 The History of Parallel Operating System	3
1.2 Multiprocessor System Architecture	6
1.3 Classify of Parallel Operating System	9
1.3.1 Control Method	9
1.3.2 Kernel Structure	10
Chapter 2 Process and Thread	14
2.1 Process Description	14
2.1.1 What's the Process	15
2.1.2 Operating System Control Structure	15
2.1.3 Process Control Structure	17
2.2 Process State	19
2.2.1 Create and Destroy of Process	20
2.2.2 Change Model of Process State	21
2.2.3 Process Suspend	23
2.3 Process Control	25
2.3.1 Execute Model	25
2.3.2 Process Switch	26
2.3.3 Execute Model of Operating System	27
2.4 Multiprocessor and Thread	30
2.4.1 Symmetry Multiprocessor	30
2.4.2 Process Model and Thread Model	32
2.4.3 Thread Implementation	40
Chapter 3 Processor Schedule	45
3.1 Process Schedule	45
3.1.1 Scheduling by Level	46
3.1.2 Schedule Method and Implementation	47
3.1.3 Schedule Arithmetic	49
3.2 Schedule Design on Multiprocessor	53
3.2.1 Schedule and Structure of Parallel System	54
3.2.2 Thread Schedule	57
3.3 Allocation and Schedule on Distributed Memory Multiprocessor	62

3.3.1	Classify of Distributed Memory Multiprocessor	62
3.3.2	Summary of Processor Allocation and Schedule Arithmetic	63
3.3.3	Typical Example of Allocation and Schedule Arithmetic	66
3.4	Typical Analyze of Schedule Parallel Operating System Processor	70
3.4.1	LWP Pure Kernel Thread and User Level Thread	70
3.4.2	Processor Schedule in UNIX SVR4.2MP	71
3.4.3	Processor in Schedule Windows NT	73
Chapter 4	Memory Management	75
4.1	Memory Access Models	75
4.1.1	Uniform Memory Access	75
4.1.2	Non-Uniform Memory Access	76
4.1.3	Non-Remote Memory Access	77
4.1.4	Cache Only Memory Access	77
4.2	Virtual Memory	78
4.2.1	Address Space	78
4.2.2	Address Translation	79
4.2.3	UMA Model Memory Management	81
4.2.4	The Others Model Memory Management	81
4.3	Data Coherence Management	82
4.3.1	Coherence Model	82
4.4	Memory Coherence Model and Programming Model	88
4.4.1	Cache Conherence Protocol	89
4.5	Shared Virtual Memory System and Single Memory Space Programming Model	90
4.6	Data Accessing and Localizing in Parallel Environment	92
4.7	Page Migration Technology	92
4.7.1	Four Key Problem About Page Migration	93
4.7.2	Page Fault Trigger Page Migration Technology	95
4.7.3	Cache Invalidation Trigger Page Migration Technology	97
4.7.4	Lazy Page Migration	100
4.7.5	Methods Compare of Page Migration Technology	103
Chapter 5	Synchronization and Mutex	105
5.1	Origin of Synchronization	105
5.1.1	Typical Models in Single Processor System	105
5.2	Main Methods in Single Processor System	106
5.3	Particularity of Synchoronization and Mutex in Parallel Environment	109
5.3.1	Data Protect Problem	111
5.3.2	Modify of Lock Semantic	111
5.3.3	Problem About Interrupt Mask	112
5.3.4	Semaphore Limitation	112

5.3.5 Wake Up Lose Problem	113
5.3.6 Exceeding Problem	113
5.3.7 Share and Mutex Accessing	113
5.4 Synchronization and mutex of Shared Memory Multiprocessor	113
5.4.1 Basic Operating in Shared Memory Multiprocessor System Environment	114
5.4.2 Synchronization and Mutex Operating in Operating System Environment	116
5.4.3 Several Problem in Multiprocessor System	117
5.5 Synchronization and Mutex in Distributed Memory Environment	119
5.5.1 Centralized Control Method	119
5.5.2 None Interceder Compete Method	119
5.5.3 Token-Ring Method	120
5.5.4 Solve Method	120
5.6 Some Examples	120
5.6.1 UNIX System V Release 4	121
5.6.2 Others Implementation	122
Chapter 6 Inter Processor Communication Technology	126
6.1 Basic Concept	126
6.1.1 Message Passing	126
6.1.2 Share Variable	127
6.2 Communication Protocol	128
6.2.1 System Architecture of Layered Network	128
6.2.2 ISO/OSI Seven Level Communication Reference Model	132
6.2.3 TCP/IP Network Reference Model	133
6.3 MPP Operating System Communication Technology	135
6.3.1 PVM Parallel Virtua System	135
6.3.2 Message Passing Interface Standard	139
6.4 Active Message Communication and Application Programming Interface	142
6.4.1 Brief Introduction	142
6.4.2 Active Message API: End and Group	145
6.4.3 Active Message API: Parallel and Synchronization	150
6.4.4 Active Message API: Transfer Operating	154
6.4.5 Solaris/Myrinet End Administration	157
6.5 Multiple Communication Methods Technology and NEXUS Implementation	158
6.5.1 Reason	158
6.5.2 Multiple Methods Communication	159
6.5.3 NEXUS Implementation	159
Chapter 7 IO Parallel File System and Parallel I/O	164
7.1 Distributed File System	165
7.1.1 Distributed File System Design	166

7.1.2	Distributed File System Implementation	168
7.1.3	New Technology on Distributed File System	172
7.1.4	Distributed File System Example	173
7.2	Parallel I/O and Parallel File System	181
7.2.1	File Abstract	181
7.2.2	Parallel I/O Method	183
7.2.3	Parallel File System	184
7.2.4	Galley Parallel File System	184
7.2.5	PVFS Parallel File System	191
7.2.6	Panda Parallel I/O Library	194
7.2.7	MPI-IO	196
Chapter 8	Mach Operating System Kernel Analysis	200
8.1	Mach Summary	200
8.2	Mach Micro kernel	201
8.2.1	Micro Kernel	201
8.2.2	Object in Mach Kernel	202
8.2.3	C/S Model	202
8.3	Process Management of Mach	203
8.3.1	Task	203
8.3.2	Thread	203
8.3.3	Schedule of Mach Thread	205
8.4	Memory Management of Mach	206
8.4.1	Memory Object	207
8.4.2	Virtual Memory Object of Mach	207
8.4.3	Mapping of Memory Object to Virtual Space	208
8.5	Communication System of Mach	209
8.5.1	Basic Concept	209
8.5.2	Implementation of Port Right Protection	209
8.5.3	Generic Method of Port Right Protection	212
8.5.4	Network Communication	213
References	214