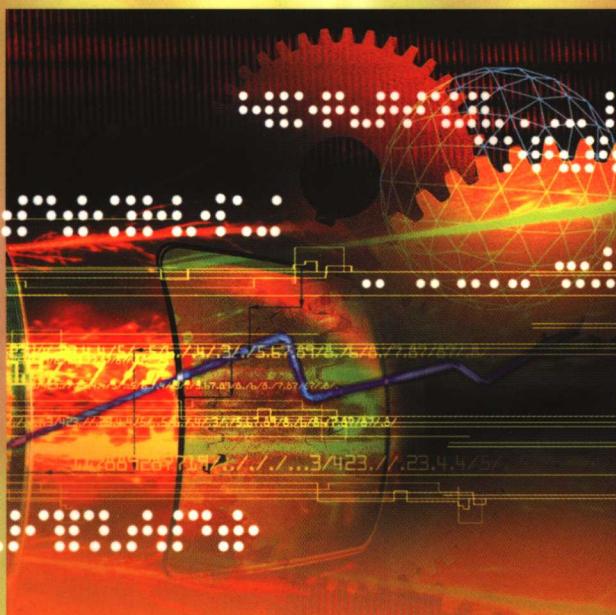


● 高等学校研究生系列教材

# 分布计算系统

## Distributed Computing Systems

徐高潮 胡亮 鞠九滨



高等教育出版社  
HIGHER EDUCATION PRESS

高等学校研究生系列教材

# 分布计算系统

徐高潮 胡 亮 鞠九滨

高等教育出版社

## 内 容 提 要

本书介绍由计算机网络构成的分布计算系统的结构和实现技术,侧重于基本概念、基本原理和基本方法的讲授。全书共分为12章:分布计算系统的基本概念和体系结构,分布计算系统的进程通信,分布式程序设计语言,命名与保护,分布式同步和互斥机构,死锁问题及其处理技术,容错技术,分布式数据管理,分布式文件系统的设计问题与实现方法,分布式调度,分布式共享存储器技术以及基于对象的分布式系统。

本书比较全面和系统地反映了国际上在这一领域取得的主要成果。特别是最近几年的成果,强调了通信、容错、死锁、同步和互斥、并发控制、分布式调度及分布式对象方面的内容。

本书叙述详细、由浅入深,具有丰富的图示,力求做到既通俗易懂,又具有一定的理论深度。本书可作为高等院校本科高年级学生和研究生教材,也可供有关科技人员参考。

## 图书在版编目(CIP)数据

分布计算系统/徐高潮,胡亮,鞠九滨.一北京:  
高等教育出版社,2004.1

ISBN 7-04-013309-1

I. 分... II. ①徐... ②胡... ③鞠... III. 分布式  
计算机系统 IV. TP338.8

中国版本图书馆CIP数据核字(2003)第113171号

---

出版发行 高等教育出版社  
社 址 北京市西城区德外大街4号  
邮政编码 100011  
总 机 010-82028899

购书热线 010-64054588  
免费咨询 800-810-0598  
网 址 <http://www.hep.edu.cn>  
<http://www.hep.com.cn>

经 销 新华书店北京发行所  
印 刷 北京二二〇七工厂

开 本 787×1092 1/16 版 次 2004年1月第1版  
印 张 27.25 印 次 2004年1月第1次印刷  
字 数 530 000 定 价 33.80元

---

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换。

**版权所有 侵权必究**

## 前　　言

由于计算机技术和网络技术的快速发展以及分布计算系统的应用领域日益扩大,人们对分布计算系统的研究兴趣日益增加,现在,分布计算系统已经是一个非常热门的研究领域。在过去的 10 年中,分布计算系统的理论体系和相关技术更加成熟,许多研究者在此领域进行了大量的研究工作,并取得了大量的研究成果。在此期间,作者在教学过程中也积累了许多新的素材,同时也积累了自己在分布计算系统方面的一些研究成果,本书力求反映这些最新的研究成果。分布式应用已涉及到社会生活的各个方面,越来越多高等院校的计算机系开设了分布计算系统这门课程。

本书所选用的参考文献主要有:《Distributed Systems: Principles and Paradigms》(Andrew S. Tanenbaum 著,清华大学出版社 2002 年影印版)、《分布式系统设计》(吴杰著,机械工业出版社 2001 年中译本)、《Distributed Operating Systems: Concepts and Practice》(Doreen L. Galli 著,人民邮电出版社影印版)、《机群计算》(鞠九滨等著,吉林大学出版社 1999 年)等。为了保证本书内容的正确性、准确性和时代性,作者选用参考资料时,尽量选用了最近出版的有权威性的原著和最新的研究成果,包括文献和专著。在此,我们对所引用著作的作者和出版者深表感谢。

本书介绍了用计算机网络组成的分布计算系统的结构和实现技术,侧重于基本概念、基本原理和基本方法。全书共分为 12 章:分布计算系统的概念和体系结构、分布计算系统的进程通信、分布式程序设计语言、命名与保护、分布式同步和互斥机构、死锁问题及其处理技术、容错技术、分布式数据管理、分布式文件系统的设计问题与实现方法、分布式调度、分布式共享存储器技术以及基于对象的分布式系统。

读者应该在学过计算机网络、计算机组成原理和操作系统之后学习本课程。如果读者已经学过网络程序设计课程,则 2.1 节和 2.2 节可以略去不学。目录中带有 \* 标记的部分可选择讲授。

清华大学郑纬民教授审阅了全稿,对本书提出了不少宝贵意见,作者在此表示诚挚的谢意。

衷心希望读者指出错误和提出改进意见。

作者

吉林大学计算机科学与技术学院

2003 年 9 月

**策划编辑** 刘建元  
**责任编辑** 武林晓  
**封面设计** 王凌波  
**责任印制** 宋克学

## 郑重声明

高等教育出版社依法对本书享有专有出版权。任何未经许可的复制、销售行为均违反《中华人民共和国著作权法》，其行为人将承担相应的民事责任和行政责任，构成犯罪的，将被依法追究刑事责任。为了维护市场秩序，保护读者的合法权益，避免读者误用盗版书造成不良后果，我社将配合行政执法部门和司法机关对违法犯罪的单位和个人给予严厉打击。社会各界人士如发现上述侵权行为，希望及时举报，本社将奖励举报有功人员。

**反盗版举报电话：**(010) 58581897/58581698/58581879/58581877

**传 真：**(010) 82086060

**E - mail:** dd@hep.com.cn 或 chenrong@hep.com.cn

**通信地址：**北京市西城区德外大街 4 号

高等教育出版社法律事务部

**邮 编：**100011

**购书请拨打电话：**(010)64014089 64054601 64054588

# 目 录

<b>第一章 绪论</b> .....	( 1 )
1.1 为什么需要分布计算系统	..... ( 1 )
1.2 分布计算系统的相关概念	..... ( 2 )
1.2.1 什么是分布计算系统	..... ( 2 )
1.2.2 松散耦合与紧密耦合分布计算 系统	..... ( 3 )
1.2.3 同构型与异构型分布计算 系统	..... ( 4 )
1.3 分布计算系统的优点和新问题	..... ( 6 )
1.3.1 分布计算系统的优点	..... ( 6 )
1.3.2 分布计算系统的新问题	..... ( 7 )
1.4 分布计算系统的透明性	..... ( 7 )
1.4.1 透明性的概念	..... ( 7 )
1.4.2 影响透明性的因素	..... ( 8 )
1.5 分布计算系统与计算机网络系统	... ( 10 )
1.5.1 网络操作系统与分布式操作 系统	..... ( 10 )
1.5.2 计算机网络系统与分布计算 系统的区别	..... ( 12 )
1.6 分布计算系统的体系结构与设计 问题	..... ( 14 )
1.6.1 分布计算系统的分层体系 结构	..... ( 14 )
1.6.2 分布计算系统的组成	..... ( 16 )
1.6.3 基于中间件的分布计算系统	..... ( 17 )
1.6.4 分布计算系统的设计问题	..... ( 19 )
习题	..... ( 22 )
参考文献	..... ( 22 )
<b>第二章 进程通信</b> .....	( 23 )
2.1 同一节点上的进程间通信	..... ( 23 )
2.1.1 管道	..... ( 23 )
2.1.2 消息队列	..... ( 26 )
2.1.3 共享内存	..... ( 30 )
2.2 不同节点上的进程间通信	..... ( 31 )
2.2.1 网络通信分层结构模型	..... ( 31 )
2.2.2 进程通信原语	..... ( 36 )
2.2.3 报文传递实例 1:socket 进程 通信	..... ( 41 )
2.2.4 报文传递实例 2:MPI 进程 通信	..... ( 46 )
2.2.5 RPC 实例 1:SUN RPC	..... ( 48 )
2.2.6 RPC 实例 2:DCE RPC	..... ( 51 )
2.3 组通信	..... ( 54 )
2.3.1 组通信的概念	..... ( 54 )
2.3.2 组通信的设计问题	..... ( 56 )
2.3.3 ISIS 中的组通信	..... ( 59 )
习题	..... ( 62 )
参考文献	..... ( 63 )
<b>第三章 分布式程序设计语言</b> .....	( 64 )
3.1 分布式程序设计语言概述	..... ( 64 )
3.1.1 分布式应用程序的分类	..... ( 64 )
3.1.2 分布式程序设计与顺序程序 设计的区别	..... ( 65 )
3.1.3 分布式程序设计语言的分类	..... ( 66 )
3.2 并行性的支持	..... ( 68 )
3.2.1 并行性的概念	..... ( 68 )
3.2.2 并行性的表示	..... ( 69 )
3.2.3 并行计算到物理处理机的 变换	..... ( 72 )
3.3 进程通信与同步的支持	..... ( 74 )

3.3.1 报文传递.....	(74)	4.2.2 公开密钥加密方法 .....	(127)
3.3.2 共享数据.....	(78)	4.3 保护 .....	(129)
3.3.3 非确定性的表示和控制.....	(80)	4.3.1 保护的目标与要求 .....	(129)
3.4* 逻辑上分布地址空间的语言 .....	(83)	4.3.2 公开密钥加密技术实现数字 签名 .....	(131)
3.4.1 同步式报文传递语言.....	(83)	4.3.3 单密钥加密技术实现数字 签名 .....	(132)
3.4.2 异步式报文传递语言.....	(84)	4.3.4 使用报文摘要的数字签名 .....	(133)
3.4.3 基于会合的语言.....	(85)	4.3.5 权能的保护 .....	(133)
3.4.4 基于远程过程调用的语言.....	(86)	4.3.6 分布系统中访问位置的控制...	(136)
3.4.5 多重通信原语.....	(87)	4.4 保护的例子:Amoeba .....	(137)
3.4.6 基于对象的语言.....	(88)	4.4.1 信口 .....	(137)
3.4.7 基于原子事务处理的语言.....	(89)	4.4.2 权能 .....	(139)
3.5* 逻辑上共享地址空间的语言 .....	(91)	4.4.3 用软件 F 盒保护 .....	(140)
3.5.1 并行函数式语言.....	(91)	习题 .....	(141)
3.5.2 并行逻辑语言.....	(92)	参考文献 .....	(142)
3.5.3 基于分布数据结构的语言.....	(93)	<b>第五章 同步和互斥</b> .....	(143)
3.6 分布式控制描述语言 DCDL .....	(96)	5.1* 分布式系统中的资源管理 .....	(143)
3.6.1 DCDL 中并行性的表示 .....	(96)	5.1.1 资源管理方式 .....	(143)
3.6.2 选择语句.....	(97)	5.1.2 控制空间 .....	(144)
3.6.3 重复语句.....	(97)	5.1.3 分散控制与通信 .....	(148)
3.6.4 语句并发(或并行)的条件 .....	(99)	5.1.4 资源的分配原则 .....	(149)
3.6.5 DCDL 中的通信 .....	(99)	5.2 同步机构 .....	(149)
3.6.6 DCDL 中的通信容错 .....	(101)	5.2.1 分布式系统中同步机构的作 用 .....	(149)
习题 .....	(102)	5.2.2 分布计算系统中的同步机构...	(151)
考文献 .....	(103)	5.2.3* 物理时钟 .....	(152)
<b>第四章 命名与保护</b> .....	(105)	5.2.4 逻辑时钟 .....	(156)
4.1 分布式系统中的命名 .....	(105)	5.3 系统的全局状态 .....	(161)
4.1.1 名字、标识符和地址 .....	(105)	5.3.1 全局状态的形式定义 .....	(162)
4.1.2 分布式系统中的名字 .....	(107)	5.3.2 全局状态的获取 .....	(163)
4.1.3 名字的结构 .....	(108)	5.3.3 一致全局状态的充要条件 .....	(164)
4.1.4 名字空间 .....	(109)	5.4 互斥算法 .....	(165)
4.1.5 名字解析 .....	(111)	5.4.1 互斥问题 .....	(165)
4.1.6 分布式系统中的名字空间 的实现 .....	(114)	5.4.2 集中式互斥算法 .....	(166)
4.1.7 实例:DNS .....	(119)	5.4.3 非基于令牌的互斥算法 .....	(167)
4.2 加密技术 .....	(123)		
4.2.1 传统加密方法 .....	(124)		

---

5.4.4 基于令牌的互斥算法 .....	(171)	7.4.4 混合检查点 .....	(214)
5.4.5* 选举算法 .....	(174)	7.4.5 报文日志 .....	(214)
5.4.6* 自稳定算法 .....	(176)	7.5 拜占庭故障的恢复 .....	(216)
习题 .....	(178)	7.5.1 恢复中的设计问题 .....	(216)
参考文献 .....	(179)	7.5.2 错误屏蔽和进程复制 .....	(218)
<b>第六章 分布式系统中的死锁 .....</b>	<b>(182)</b>	7.5.3 容错系统中的一致性协议 .....	(219)
6.1 死锁问题 .....	(182)	7.6 可靠的组通信 .....	(225)
6.1.1 死锁发生的条件 .....	(182)	7.6.1 基本的可靠组播方案 .....	(225)
6.1.2 死锁的图论模型 .....	(183)	7.6.2 可靠的组播通信中的可扩	
6.1.3 处理死锁的策略 .....	(184)	充性 .....	(227)
6.1.4 死锁的 AND 条件和 OR 条件 .....	(185)	7.6.3 原子组播 .....	(229)
6.2 死锁的预防 .....	(186)	习题 .....	(234)
6.2.1 预防死锁的一般方法 .....	(186)	参考文献 .....	(235)
6.2.2 基于时间戳的预防死锁方法 .....	(187)	<b>第八章 分布式数据管理 .....</b>	<b>(238)</b>
6.3 死锁的检测 .....	(188)	8.1 一致性模型 .....	(238)
6.3.1 集中式死锁检测 .....	(188)	8.1.1 严格一致性 .....	(238)
6.3.2 分布式死锁检测 .....	(190)	8.1.2 顺序一致性和可线性化	
6.3.3 层级式死锁检测 .....	(191)	一致性 .....	(239)
6.3.4 死锁检测的实例 .....	(192)	8.1.3 相关一致性 .....	(241)
习题 .....	(197)	8.1.4 FIFO 一致性 .....	(242)
参考文献 .....	(198)	8.1.5 弱一致性 .....	(244)
<b>第七章 分布式系统中容错技术 .....</b>	<b>(200)</b>	8.1.6 释放一致性 .....	(245)
7.1 分布式系统中的故障模型 .....	(200)	8.1.7 进入一致性 .....	(247)
7.1.1 基本概念 .....	(200)	8.2 并发控制 .....	(249)
7.1.2 基本的故障模型 .....	(201)	8.2.1 并发控制的目标与事务处理 .....	(249)
7.2 容错系统的基本构件 .....	(204)	8.2.2 可串行化调度 .....	(253)
7.2.1 坚固存储器 .....	(204)	8.2.3 基于锁的并发控制 .....	(257)
7.2.2 故障 - 停止处理器 .....	(204)	8.2.4 基于时间戳的并发控制 .....	(259)
7.2.3 原子操作 .....	(205)	8.2.5 乐观的并发控制 .....	(261)
7.3 节点故障的处理 .....	(205)	8.3 原子事务处理 .....	(261)
7.3.1 向后式恢复 .....	(206)	8.3.1 原子事务处理的性质 .....	(261)
7.3.2 向前式恢复 .....	(208)	8.3.2 事务处理的分类 .....	(263)
7.4 检查点算法 .....	(209)	8.3.3 原子事务处理的实现 .....	(264)
7.4.1 一致性检查点 .....	(209)	8.3.4 基于原子事务处理的局部	
7.4.2 异步检查点 .....	(211)	恢复 .....	(264)
7.4.3 同步检查点 .....	(212)	8.3.5 分布式提交协议 .....	(267)

---

8.4 多副本更新和一致性管理 .....	(269)	9.7.5 NFS 中的文件封锁 .....	(308)
8.4.1 分布式系统中的系统数据库…	(270)	9.7.6 缓存和复制 .....	(310)
8.4.2 兼容可串行化 .....	(271)	9.7.7 NFS 中的容错 .....	(311)
8.4.3 主站点方法 .....	(272)	9.7.8 NFS 的安全性 .....	(313)
8.4.4 循环令牌方法 .....	(273)	9.8* 其他的分布式文件系统及 其比较 .....	(316)
8.4.5 同步表决方法 .....	(273)	9.8.1 设计目标 .....	(316)
8.4.6* 活动复制控制方法 .....	(275)	9.8.2 通信和进程 .....	(317)
8.4.7 法定数方法 .....	(276)	9.8.3 命名 .....	(318)
习题 .....	(279)	9.8.4 同步 .....	(318)
参考文献 .....	(279)	9.8.5 缓存和复制 .....	(319)
<b>第九章 分布式文件系统 .....</b>	<b>(282)</b>	9.8.6 容错 .....	(319)
9.1 分布式文件系统的特点和基本 要求 .....	(282)	9.8.7 安全性 .....	(320)
9.1.1 分布式文件系统的特点 .....	(282)	习题 .....	(321)
9.1.2 分布式文件系统的基本要求…	(283)	参考文献 .....	(321)
9.2 分布式文件系统中的命名 .....	(284)	<b>第十章 分布式调度 .....</b>	<b>(324)</b>
9.2.1 命名方案 .....	(285)	10.1 调度算法概述 .....	(324)
9.2.2 命名的实现技术 .....	(286)	10.1.1 调度算法的分类 .....	(324)
9.3 共享语义 .....	(288)	10.1.2 调度算法的目标和有效性 评价 .....	(325)
9.4 缓存 .....	(290)	10.2 静态调度 .....	(327)
9.4.1 文件的远程访问方法 .....	(290)	10.2.1 任务划分与分配 .....	(327)
9.4.2 缓存的粒度和地点 .....	(291)	10.2.2 基于任务优先图的任务 调度 .....	(331)
9.4.3 更新策略、缓存有效性和 一致性 .....	(292)	10.2.3 两种最优调度算法 .....	(333)
9.4.4 缓存和远程服务的比较 .....	(293)	10.2.4 基于任务相互关系图的 任务调度 .....	(335)
9.5 容错和可扩充性 .....	(294)	10.3 动态调度 .....	(338)
9.5.1 无状态服务和有状态服务 .....	(294)	10.3.1 动态调度的组成要素 .....	(338)
9.5.2 可用性与文件复制 .....	(296)	10.3.2 动态负载平衡算法的分类、设计 决策和使用的参数 .....	(340)
9.5.3 可扩充性 .....	(297)	10.4 空闲工作站的调度结构 .....	(343)
9.5.4 用线程实现文件服务员 .....	(298)	10.4.1 工作站共享问题 .....	(343)
9.6 安全性 .....	(299)	10.4.2 工作环境 .....	(345)
9.7 SUN 网络文件系统(NFS) .....	(299)	10.4.3 集中式调度 .....	(346)
9.7.1 NFS 概述 .....	(300)	10.4.4 分散式调度 .....	(348)
9.7.2 NFS 中的通信 .....	(303)		
9.7.3 NFS 服务员 .....	(304)		
9.7.4 NFS 中的命名 .....	(305)		

---

10.4.5 混合式调度	(349)	11.3.2 全映像目录	(375)
10.5 进程转移和远程执行	(350)	11.3.3 有限目录	(377)
10.5.1 进程转移和远程执行的目的和方法	(350)	11.3.4 链式目录	(378)
10.5.2 Sprite 的进程迁移和远程执行设备	(351)	11.3.5 只对专用数据进行缓存的方案	(379)
10.5.3 V 系统中的可抢先的远程执行设备	(354)	11.3.6 性能比较	(379)
10.5.4 NEST 中的透明的远程执行设备	(354)	11.4* DSM 系统的实现	(380)
10.6* 空闲工作站共享系统 Sidle	(355)	11.4.1 实现 DSM 的基本方法	(380)
10.6.1 Sidle 的组成	(355)	11.4.2 结构和粒度	(381)
10.6.2 Sidle 的调度	(356)	11.4.3 数据定位和访问	(382)
10.6.3 Sidle 的透明远程执行设备	(357)	11.4.4 一致性协议	(383)
习题	(359)	11.4.5 替换策略	(385)
参考文献	(359)	11.4.6 颠簸	(386)
<b>第十一章 分布式共享存储器</b>	(362)	11.4.7 可扩充性	(386)
11.1 基本概念	(362)	11.4.8 异构性	(387)
11.1.1 什么是分布式共享存储器系统	(362)	11.4.9 其他有关算法	(387)
11.1.2 为什么需要分布式共享存储器	(363)	11.5* DSM 实例:Ivy 和 MemNet	(388)
11.1.3 共享存储器中缓存一致性方法	(364)	11.5.1 Ivy——软件实现的 DSM	(388)
11.1.4 DSM 的设计与实现问题	(365)	11.5.2 Ivy 一致性协议	(388)
11.1.5 一致性语义	(366)	11.5.3 Ivy 存储器管理	(391)
11.2 实现 DSM 的算法	(367)	11.5.4 Ivy 中的进程同步	(392)
11.2.1 算法使用的模型和环境	(367)	11.5.5 MemNet——硬件实现的 DSM	(392)
11.2.2 中央服务员算法	(368)	11.5.6 MemNet 缓存一致性协议	(393)
11.2.3 迁移算法	(369)	11.5.7 Ivy 与 MemNet 的比较	(394)
11.2.4 读复制算法	(370)	习题	(395)
11.2.5 全复制算法	(371)	参考文献	(395)
11.2.6 算法性能	(372)	<b>第十二章 基于对象的分布式系统</b>	(397)
11.2.7 算法比较	(373)	12.1 分布式对象	(397)
11.3 使用目录的 DSM	(374)	12.1.1 对象的概念	(397)
11.3.1 目录方案的分类	(375)	12.1.2 对象的类型	(399)
		12.2 CORBA	(400)
		12.2.1 CORBA 的总体结构	(400)
		12.2.2 CORBA 的对象模型	(401)
		12.2.3 接口库和实现库	(402)
		12.2.4 CORBA 的服务	(403)

12.2.5 CORBA 的通信 .....	(404)	12.3.6 DCOM 的 Moniker .....	(418)
12.2.6 CORBA 的 POA .....	(409)	12.4* Clouds 系统 .....	(419)
12.3 DCOM .....	(411)	12.4.1 Clouds 的对象 .....	(419)
12.3.1 COM 和 DCOM .....	(411)	12.4.2 Clouds 的线程 .....	(420)
12.3.2 DCOM 的对象模型 .....	(412)	12.4.3 Clouds 的存储器 .....	(420)
12.3.3 DCOM 的类型库和注册 .....	(413)	习题 .....	(421)
12.3.4 DCOM 的服务 .....	(415)	参考文献 .....	(421)
12.3.5 DCOM 的通信 .....	(415)		

# 第一章 绪 论

## 1.1 为什么需要分布计算系统

20世纪90年代以来,两大技术的快速进步导致了人们对分布计算系统的兴趣迅速增加,使得分布计算系统成为一个热点研究领域。这两大技术进步分别是:第一,计算机硬件技术和软件技术的发展;第二,高性能计算机网络技术的发展。这两大技术的进步不仅极大地改变了人们使用计算机的方式,同时使得人们对尽可能多的计算能力、数据的透明访问、高性能和高可靠性的目标的追求不再是梦想。

计算机技术的发展可以通过人们使用计算机的方式的改变来描述。在20世纪50年代,程序设计人员使用计算机前必须预约上机时间,他们使用计算机时,会占用全部计算机资源。60年代出现了批处理技术,人们提交作业并排队等候处理,计算机每次运行一个作业,用户晚些时候来取结果。70年代产生了分时系统,人们可同时使用一台计算机,而每个用户都觉得他在单独使用整个计算机。80年代是个人计算机的时代,每个用户有一台属于自己专用的计算机。现在,不仅每个用户有自己专用的计算机,而且通过计算机网络可同时使用多台计算机。

使用计算机的方式由预约上机方式转变为批处理方式的首要技术条件是必须有足够的大的计算机存储器。只有存储器做得充分大,才可以装得下整个操作系统以及待处理的应用程序。在批处理方式中,缓冲、假脱机和多道程序的脱机处理等技术得到了广泛应用。需要批处理技术的原因是它可以有效地利用昂贵的计算机机时,减少了计算机的空闲时间。

出现分时系统的主要技术条件是计算机变得更为便宜而且功能更为强大。使用分时系统的原因是因为它可以使程序设计人员获得较高的效率,也更进一步地提高了计算机的利用率。分时系统是迈向分布计算系统的一步,用户可以在不同的地点共享和访问资源,但分时系统给用户提供的是一个集中的共享环境,允许很多设备如打印机、存储器以及软件和数据共享。

超大规模集成电路和高性能计算机网络技术的发展使得个人计算机代替了分时系统。由于这时每个用户都有自己专用的计算机,并且可通过计算机网络共享和访问其他的计算机资源,使得用户在任何时间都能保证使用所需的计算能力成为可能,但此时提供给用户的是

一个基于网络的、分散的共享环境。

与分时系统的共享环境相比,在基于网络的分散的共享环境中,用户在系统的使用和管理上会遇到新的问题。在使用上,个人计算机的操作系统软件允许在网络上的个人计算机之间复制文件和远程登录,要求用户必须知道本地对象和远程对象之间的差别,以及对象放在哪个远程机器上。如果网络的规模很大,这个问题将变得很严重。

在系统管理上,分时系统的操作人员每天进行文件系统备份操作,系统管理员把可用的处理机分配给最需要处理机的用户,系统编程人员只需简单地安装新的或已改进的软件。但在由个人计算机互连的网络环境中,个人计算机的用户必须既是一个操作员,又是系统管理员和系统编程人员,在拥有成百上千台个人计算机的建筑物中,操作人员不可能四处奔走进行文件备份操作。系统编程人员也不再能很简单地把新的软件放到文件系统中。

针对上述使用和管理上出现的新问题,虽然也有一系列的解决方法,但没有一个能像使用分时系统共享环境那样令人满意。例如,最常用的一种方法是在个人计算机系统中增加一个网络复制命令,使用此命令把文件从一个计算机复制到另一个计算机上。另一种稍好的方法是使用网络文件系统,它允许文件在某种程度上实现真正共享。但在所有的方法中,除极少数以外,用户总要知道本地操作和远程操作的差别。出现这些问题的最主要原因是,传统的操作系统(它仍然是现代个人计算机软件的基础)从来都不是为具有多个处理机和很多文件系统的环境所设计的,这种环境需要分布式操作系统。

在分布计算系统中,多台计算机构成一个完整的系统,其行为类似一个单机系统。即登录到系统的用户不必了解系统中有多少台机器,它们的位置在哪里,它们的功能是什么,文件在哪里,作业在哪一台机器上运行等任何有关硬件物理分布的细节。分布式操作系统是实现分布计算系统的核心。

以上是从人们使用计算机方式的变化这个角度上分析了为什么需要分布计算系统。其实人们需要分布计算系统还在于它和单机系统相比具有许多优点。

## 1.2 分布计算系统的相关概念

### 1.2.1 什么是分布计算系统

分布计算系统又叫分布式计算机系统,简称分布式系统。对于分布计算系统,不同的研究者给出了不同的定义,但没有一个是完全令人满意的,也没有一个被所有的研究者所承认的。例如,Andrew S. Tanenbaum 教授给出的定义如下:分布计算系统是由多个独立的计算机系统相互连接而成的计算系统,从用户的角度来看,它好像是一个集中的单机系统[7]。

总结不同研究者给出的不同定义,可以概括如下:

分布计算系统是由多个相互连接的处理资源组成的计算系统,它们在整个系统的控制下可合作执行一个共同的任务,最少依赖于集中的程序、数据和硬件。这些处理资源可以是物理上相邻的,也可以是在地理上分散的。

现在对这一定义进行说明:

- (1) 系统是由多个处理器或计算机系统组成的。
- (2) 这些计算资源可以是物理上相邻的、由机器内部总线或开关连接的处理器,通过共享主存进行通信;这些计算资源也可以是在地理上分开的、由计算机通信网络(远程网或局域网)连接的计算机系统,使用报文(Message)进行通信。
- (3) 这些资源组成一个整体,对用户是透明的,即用户使用任何资源时都不必知道这些资源在哪里。
- (4) 一个程序可分散到各个计算资源上运行。
- (5) 各个计算机系统地位平等,除了受全系统的操作系统控制外,不存在主从控制和集中控制环节。

分布计算系统属于多指令流多数据流(MIMD)结构。

### 1.2.2 松散耦合与紧密耦合分布计算系统

从结构上分,有两大类分布式系统:紧密耦合分布式系统和松散耦合分布式系统。紧密耦合分布式系统由多个处理机经机器内部总线或机器内互连网络(如 mesh、cube 等)连接而成,因此也叫多处理机系统。紧密耦合分布式系统中的多个处理机有共享的主存储器(如图 1.2.1(a)和 1.2.1(b)),也可以既有共享的主存储器又有自己专用的主存储器(如图 1.2.1(c))。紧密耦合分布式系统一般局限于一个局部区域,各部分相距很近,可以说在物理上是分散的,但不是在地理上分散的。各处理机通过共享主存交换信息,并具有较高的通信速

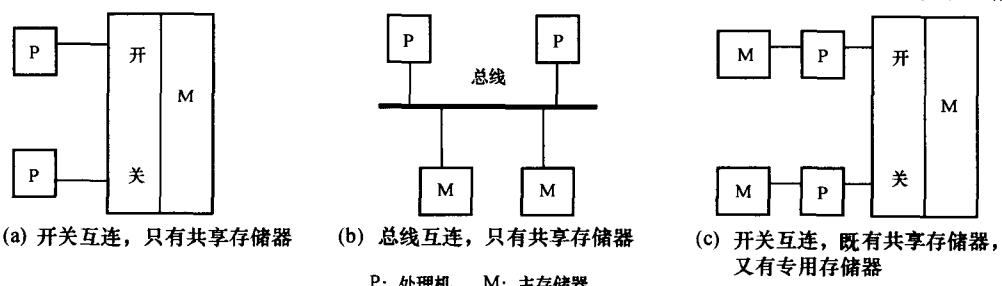


图 1.2.1 紧密耦合分布式系统

度。本书不讨论这种系统。

松散耦合分布式系统由多个计算机系统(包括处理机和相应的主存储器)经过通信网络连接而成,是多计算机系统(如图 1.2.2)。计算机之间或处理机之间使用报文交换方法通信。松散耦合分布式系统没有专门用于各处理机共享的主存,但是,各计算机的主存可以共享,构成所谓分布式共享存储器系统。通信网络可以是远程网,也可以是局域网。通常的“分布式系统”一词指的是松散耦合分布式系统,本书以后各部分的分布式系统除特别说明外,均指松散耦合分布式系统。前面 Andrew S. Tanenbaum 教授对分布计算系统的定义,其实就是指松散耦合分布式系统。

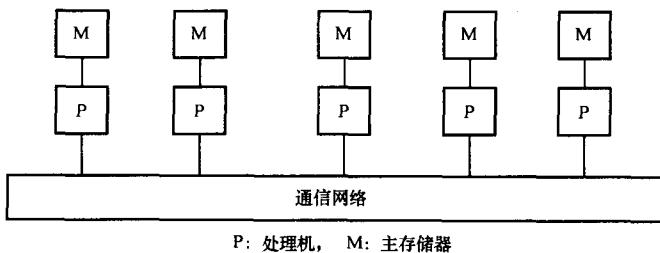


图 1.2.2 松散耦合分布式系统

### 1.2.3 同构型与异构型分布计算系统

按照组成分布计算系统的计算机种类和所使用的计算机网络的种类,分布计算系统可分成同构型分布式系统和异构型分布式系统。对于同构型分布式系统而言,组成该系统的计算机的硬件和软件是相同的或非常相似的,同时,组成该系统的计算机网络的硬件和软件也是相同的或非常相似的。对于异构型分布式系统而言,组成该系统的计算机的硬件或软件是不同的,或者组成该系统的计算机网络的硬件或软件也是不同的[3]。

分布计算系统的异构性包括计算机系统的异构性和通信网络的异构性。计算机系统的异构性主要有三个方面的表现,即:硬件异构性、操作系统异构性和程序设计语言异构性。

计算机系统的硬件异构性主要有以下三个方面的表现:

(1) 计算机的指令系统不同。这意味着一种机器上的程序模块不能在另一种不兼容的机器上执行,很显然,一种机器上的可执行代码程序也不能在另一种不兼容的机器上执行。

(2) 数据表示方法不同。例如,不同类型的计算机虽然都是按字节编址的,但是高字节和低字节的规定可能恰好相反。浮点数的表示方法也常常不一样。M6800 CPU 的字长是 16 位,并且可按字节单位寻址,它是高字节在前,低字节在后;而 Intel 80286 CPU 也是按字节单位寻址,但它是低字节在前,高字节在后,因此,在 Intel 80286 CPU 上运行的程序不能直接操

作在 M6800 CPU 上运行的程序所产生的数据。

(3) 机器的配置不同。尽管机器的类型可能相同,其硬件配置也可以互不兼容。例如,一台机器可能比另一台机器具有更多的内存或外存,而另一台则可能多配置了某种外部设备等。

即使硬件配置完全相同,如果运行的操作系统不同,则由这些机器组成的分布计算系统仍为异构型。操作系统异构性的三个主要表现方面为:

(1) 操作系统所提供的功能可能大不相同。例如,不同的操作系统提供了不同的命令集。

(2) 操作系统所提供的系统调用在语法、语义和功能方面不相同。

(3) 文件系统不同。文件系统的不同主要表现在三个方面。第一,不同的系统采用不同的文件命名方式。例如,UNIX 系统和 MS - DOS 系统的文件系统都是层次型的,但是 UNIX 用符号“/”分隔路径名的各部分,而 MS - DOS 文件系统中的文件名的第一部分是驱动器名,后跟一个“:”符号,最后的部分采取 \* \* \* \* \*. \* \* \* 的形式。这样,UNIX 程序可以产生和使用任何 MS - DOS 文件名,反过来则不行。所以如果把两个系统容纳在一个系统中共享文件,那么许多操作将会失败。第二,文件的保护方法和程度也不相同。例如,UNIX 系统把用户分成三类:文件的所有者、同组用户和其他用户,不同用户对某个文件有不同的访问权限。而 MS - DOS 则根本无文件保护措施。因此,很难要求一个文件系统在共享另一个文件系统的文件时能够按那个文件系统的保护方法对文件进行保护。第三,不同的文件系统有不同的文件模式。例如,UNIX 只提供了字节流模式,VMS 则提供多种记录结构的文件模式。

程序设计语言的异构性表现在不同的程序设计语言用不同的方法在文件中存储数据。例如,Pascal 程序以十进制字符串的形式在文件上存放整数,而 UNIX 系统中,C 语言则以二进制形式在文件中存放整数。因此,用不同语言编写的程序之间有可能不能直接共享数据文件。

通信网络的异构性在硬件方面主要表现在,不同的计算机网络和同一种计算机之间的接口硬件可能不同,连接方法和通信方法也可能不同。通信网络的异构性在软件方面主要表现在不同计算机网络的通信协议不相同。

随着计算机技术的发展和计算机系统的广泛应用,异构性是不可避免的。即使用户的初衷是建立一个同构型分布计算系统,但往往最终也会自然地演变为异构型分布计算系统。这是因为:

(1) 分布计算系统已成为资源共享的重要形式。随着分布计算系统资源的增多,其他用户也希望加入系统,共享其资源。这些新的系统往往有着与原有系统不同的硬件和软件。

(2) 由于硬件性能的提高和其价格的下降,当扩充一个分布计算系统时,人们往往会选