陈国良. 并行计算: 结构 • 算法 • 编程 [M]. 1999.

屈彬

2019年7月12日

1 并行计算机系统及其结构模型

1.1 并行计算(2019年7月2日)

1.1.1 并行计算与计算科学

定义 1.1 (应用需求) 应用需求分为三类, 计算密集(Compute-Intensive) 型应用, 如大型科学工程计算与数值模拟; 数据密集(Data-Intensive) 型应用, 如数字图书馆、数据仓库、数据挖掘和计算可视化等; 网络密集(Network-Intensive) 型应用, 如协同工作、遥控和远程医疗诊断等。

定义 1.2 (高性能计算和通信) High Performance Computing and Communication, HPCC。

定义 1.3 (加速战略计算创新) Accelerated Strategic Computing Initiative, AS-CI.

定义 1.4 (美国能源部三大高性能计算实验室) Lawrence Livermore、Los Alamos、Sandia。

1.2 并行计算系统互连(2019年7月3日)

1.2.1 系统互连

定义 1.5 (机群网络分类) 1. 系统域网络: (System Area Network, SAN), 连接短距离 (3 25m) 内的节点。

- 2. 局域网络: (Local Area Network, LAN), 连接企事业单位内 (500m 2k-m) 内的节点。
- 3. 节点内网络: 由处理器总线、局部(本地)总线、存储器总线构成。
- 4. 小型机系统接口: (Small Computer System Interface, SCSI) 由 I/O 总 线、系统总线构成。

定义 1.6 (工作站机群) (Cluster of Workstations, COW), 通过 SAN/LAN 互连的工作站组。

1.2.2 静态互连网络

定义 1.7 (对剖宽度) (Bisection Width): 将网络划分为两等分所必须剪去的最少边数。

1.2.3 动态互连网络

定义 1.8 (总线分类) 常用总线包括 PCI、VME、Multibus、SBus、Microchannel 和 IEEE Futurebus。

- 1. 本地总线: (Local Bus), CPU 板级上的总线。
- 2. 存储器总线:存储器板级上的总线,主要指内存与 CPU 互连的总线。
- 3. **数据总线**: I/O 与通信板级上的总线,例如硬盘与内存之间的总线,以及网卡。

1.2.4 标准互连网络分类

- 1. **光纤分布式数据接口** (Fiber Distributed Data Interface, FDDI), 采用双向光纤令牌环提供 100 200Mb/s 的数据传输。
- 2. **快速以太网**主流的互连网络,第三代以太网数据传输速度可达 1Gb/s。
- 3. myrinet由 Myricom 公司生产的千兆位包开关网。
- 4. 高性能并行接口(High Performance Parallel Interface, HiPPI),主要用于构筑异构计算机系统。
- 5. **异步传输模式**(Asynchronous Transfer Mode, ATM),在光纤通信基础 上建立起来的一种新的宽带综合业务数字网(B-ISDN)的交换技术。
- 6. 无线带宽 (InfiniBand)。

1.3 并行计算机体系结构(2019年7月4日)

1.3.1 并行计算机结构模型

大型并行机系统一般可分为六类机器:

- 1. 单指令多数据流: (Single Instruction Multiple-Data, SIMD)。
- 2. 并行向量处理机: (Parallel Vector Processor, PVP)。
- 3. 对称多处理机: (Symmetric Multiprocessor, SMP)。
- 4. 大规模并行处理机: (Massively Parallel Processor, MPP)。
- 5. 工作站机群: (Cluster of Workstations, COW)。
- 6. 分布共享存储多处理机: (Distributed Shared Memory, DSM)。

除 SIMD 外,其余五种皆为 MIMD。

1.3.2 并行计算机访存模型

- 1. 均匀存储访问: (Uniform Memory Access, UMA), 其特点包括: 物理存储器被所有处理器均匀共享; 所有处理器访问任何存储单元取相同的时间; 每台处理器可带私有 cache; 外围设备也可以一定形式共享。这种系统由于高度共享资源又被称为紧耦合系统 (Tightly Coupled System)。多核同构 CPU 系统的存储模型属于典型的 UMA 模型。
- 2. 非均匀存储访问: (Nonuniform Memory Access, NUMA)。
- 3. 全高速缓存存储访问: (Cache-Only Memory Access, COMA)。
- 4. 高速缓存一致性非均匀存储访问: (Coherent-Cache Nonuniform Memory Access, CC-NUMA)。
- 5. 非远程存储访问: (No-Remote Memory Access, NORMA)。

1.3.3 并行计算机存储组织

定义 1.9 (层次存储技术) 利用程序局部性,设置多级存储系统。底层的存储器通常具有容量大、速度慢、成本低的特点,高层的存储器通常具有容量小、速度快、成本高的特点。层次存储技术出现的目的在于寻找性能与成本之间的平衡。

定义 1.10 (高速缓存一致性问题) 当处理器将新数据写入高层存储器时,需保证其他层次中存储器的数据对应的位置具有相同的副本。

2 当代并行计算机系统介绍

跳过第2章,它的内容与第1章存在大量重复。

3 并行计算性能评测

3.1 并行计算机的一些基本性能指标(2019年7月5日)

3.1.1 CPU 和存储器的某些基本性能指标

名称	符号	含意	单位
机器规模	n	处理器的数目	
时钟速率	f	时钟周期长度的倒数	MHz
工作负载	W	计算操作的数目	Mflop
顺序执行时间	T_1	程序在单处理机上的运行时间	s (秒)
并行执行时间	T_n	程序在并行机上的运行时间	s (秒)
速度	$R_n = W/T_n$	每秒百万次浮点运算	Mflops
加速	$S_n = T_1/T_n$	衡量并行机有多快	
效率	$E_n = S_n/n$	衡量处理器的利用率	
峰值速度	$R_{peak} = nR'_{peak}$	所有处理器峰值速度之积	Mflops
顺序峰值速度	$R_{peak}^{'}$	单个处理器的峰值速度	Mflops
利用率	$U = R_n / R_{peak}$	可达速度与峰值速度之比	
通信延迟	t_0	传送 0-字节或单字的时间	$\mu \mathrm{s}$
渐进带宽	$r_{ m inf}$	传送长消息的通信速率	MD/s

3.1.2 通信开销测量方法

- 1. 乒 -乓方法: (Ping-Pong Scheme) 适用于一对节点之间通信开销的测量。
- 2. 热土豆法: (Hot-Potato),也称作救火队法(Fire-Brigade),适用于测量 多个节点(两个以上)以环形方式通信,数据包传递一圈的通信开销。

理论点到点通信开销t(m) 是消息长度 m (字节) 的线性函数:

$$t(m) = t_0 + m/r_{\rm inf}$$

其中 t_0 是通信启动时间 (μ s), r_{inf} 是渐进带宽 (MB/s)。

3.1.3 机器的成本、价格与性能/价格比

内容比较零碎, 跳过。

3.2 加速比性能定律(2019年7月6日)

- 3.2.1 Amdahl 定律
- 3.2.2 Gustafson 定律
- 3.2.3 Sun 和 Ni 定律

4 并行算法的设计基础

出于 PAC 比赛的需求,暂时先跳过第三章。

4.1 并行算法的基础知识

4.2 并行计算模型(2019年7月7日)

并行计算模型是硬件和软件之间的一种桥梁(大雾,什么叫作"桥梁"),使用它能够设计分析算法,在其上高级语言能被有效地编译且能够用硬件实现。(我觉得作者没有把并行计算模型的用途讲清楚)

4.2.1 PRAM 模型

即并行随机存储机器(Parallel Random Access Machine, PRAM),是一种抽象的并行计算同步模型。在这种模型中,假定存在着一个容量无限大的共享存储器;有多个同构处理器;在任何时刻各处理器均可通过共享存储单元相互交换数据以实现同步。根据对处理器读写时序的限制,又可分类为:

- 1. 不允许同时读和同时写(Exclusive=Read and Exclusive-Write)的 PRAM 模型,记作 **PRAM-EREW**。
- 2. 允许同时读不允许同时写(Concurrent-Read and Exclusive-Write)的 PRAM 模型,记作 **PRAM-CREW**。
- 3. 允许同时读和同时写(Concurrent-Read and Concurrent-Write)的 PRAM 模型,记作 **PRAM-CRCW**。

不过,允许同时写显然是不现实的,所以对 PRAM-CRCW 模型做了进一步的约定:

- 1. **CPRAM-CRCW**:公共的 PRAM-CRCW,只允许所有的处理器同时写相同的数。
- 2. **PPRAM-CRCW**: 优先的 PRAM-CRCW, 只允许最优先的处理器先写。
- 3. **APRAM-CRCW**: 任意的 PRAM-CRCW, 允许任意处理器自由写。

4.2.2 异步 PRAM 模型

记作 APRAM,它的特点是**分相**(Phase),每个处理器都有其私有存储(书里称作"局存")、私有时钟(书里称作"局部时钟")和局部程序;处理器之间通过共享全局存储器进行通信;无全局时钟,各处理器异步地独立执行各自的指令;若处理器之间存在任何时序依赖关系(书中称作"时间依赖关系"),须通过插入同步路障(Synchronization Barrier)以保持依赖源和目标的执行顺序。

APRAM 模型具有四种访存指令:

- 1. 全局读:将全局存储单元中的内容读入私有存储单元中。
- 2. 局部操作:对局存中的数执行操作,其结果是存入局存中。
- 3. 全局写:将局存单元中的内容写入全局存储单元中。
- 4. **同步**: 同步是计算中的一个逻辑点,在该点各处理器均需等待别的处理器 到达后才能继续执行其局部程序。

4.2.3 BSP 模型

整体同步计算模型(Bulk Synchronous Parallel),字面含义是"大"同步模型,该模型使用了三个属性描述:模块(Components)、选路器(Router)和同步路障器执行时间 L(中文又称"超级步")。BSP 是一种 MIMD-DM(DM: distributed memory)模型。

在一个超级步中,每个处理器执行各自的局部计算,通过选路器接收和发送消息;然后作一全局检查(同步),以确定当前超级步是否已由所有的处理器完成;然后前进到下一个超级步。

在分析 BSP 模型的性能时,假定局部操作可在一个超级步内完成。在每一超级步中,一个处理器至多发送或接收 h 条消息(称为 h-relation)。假定 s

是传输建立时间,g 定义为每秒处理器所能完成的局部计算数目与每秒选路器所能传输的数据量之比,那么传送 h 条消息的时间开销为 gh+s,若 $gh\geq 2s$,则超级步步长 $L\geq gh$ (为什么?)。硬件可设置 L 的下限(例如宽的通信带宽减小 g),而软件可设置 L 的上限(例如并行粒度越大,L 越大)。

4.2.4 logP 模型

logP 模型是一种分布式存储的,点到点通信的多处理机模型,其中通信 网络由一组参数来描述:

- 1. l(Latency)表示在网络中消息从源到目的地所遭到的延迟。
- 2. o (Overhead)表示处理器发送或接收一条消息所需的**额外开销**。
- 3. g(gap)表示处理器可连续进行消息发送或接收的最小时间间隔。
- 4. P (Processor) 表示处理器/存储器模块数。

- 5 并行算法的一般设计策略
- 6 并行算法的基本设计技术
- 7 并行算法的一般设计过程
- 8 基本通信操作
- 9 稠密矩阵运算
- 10 线性方程组的求解
- 11 快速傅里叶变换
- 12 并行程序设计基础
- 13 共享存储系统并行编程
- 14 分布存储系统并行编程
- 14.1 基于消息传递的并行编程(2019年7月10日)

所谓基于消息传递的并行编程,是指用户必须显示地通过发送和接收消息 来实现处理器之间的数据交换。根据问题分解的两种形式,消息传递并行性的 开发也分为两种模式

- 1. **域分解**:将一个大的问题区域分解成若干个较小的问题区域,每个区域 执行同样的指令,但处理的数据不同。通常采用 SPMD (Single Process Multiple Data)编程模式实现。
- 2. **函数分解**:又称功能分解,将一个大问题分解成若干个子问题,每个子问题执行不同的指令。通常采用 MPMD(Multiple Process Multiple Data)编程模式实现。

14.1.1 SPMD 并行程序

SPMD 按照节点的地位可分为两种结构

- 1. 主机/节点结构: 节点分为 master 和 slave。
- 2. 无主机结构:每个节点都是对等的。

14.1.2 MPMD 并行程序

实现子问题之间相互作用方式有两种

- 1. **数据流方式:**(Dataflow)根据子问题数据的相关性来组织并行,按**数据 驱动**(Data-Driven)的方式决定并行执行的先后次序。
- 2. 客户/服务器方式: 就是常见的 C/S 方式。

14.2 MPI 并行编程(2019年7月10日)

(Message Passing Interface)。MPI 进程是重量级、单线程的进程。

14.2.1 最基本的 MPI

MPI 标准包含 6 个基本函数

函数名	功能
MPI_INIT	启动 MPI 计算
$MPI_FINALIZE$	结束 MPI 计算
${\rm MPI_COMM_SIZE}$	确定进程数
MPI_COMM_RANK	确定自己的进程标识符
MPI_SEND	发送一条消息
MPI_RECV	接受一条消息

下面以 MPICH2 为例说明个函数的详细信息。

MPI_INIT(int *argc, char **argv)

启动 MPI 计算, argc、argv 都是 main 函数的参数,将主函数参数传递给其他节点。

MPI_FINALIZE()

结束 MPI 计算。

MPI_COMM_SIZE(comm, *size)

输入输出类型	参数名	解释
IN	comm	通信体(handle)
OUT	size	通信体 comm 中进程的个数(integer)

MPI_COMM_RANK(comm, *pid)

输入输出类型	参数名	解释
IN	comm	通信体(handle)
OUT	pid	通信体 comm 中当前进程的标识符 (integer)

MPI_SEND(*buf, count, datatype, dest, tag, comm)

输入输出类型	参数名	解释
IN	buf	待传输数据的内存首地址(可选参数)
IN	count	待传输数据的长度(integer)
IN	datatype	待传输数据的类型(handle)
IN	dest	目标进程的 pid (integer)
IN	tag	消息标识符(integer)
IN	comm	通信体(handle)

MPI_RECV(*buf, count, datatype, source, tag, comm, status)

输入输出类型	参数名	解释
OUT	buf	数据存放位置的首地址(可选参数)
IN	count	数据的长度 (integer)
IN	datatype	待传输数据的类型(handle)
IN	source	源进程的 pid(integer)
IN	tag	消息标识符(integer)
IN	comm	通信体(handle)
OUT	status	状态对象(statuse)(这是什么?)

MPI 数据类型与 C 语言数据类型对应关系如下

MPI (C Binding)	С
MPI_BYTE	
MPI_CHAR	signed char
MPI_DOUBLE	double
MPI_INT	int
MPI_LONG	long
MPI_LONG_DOUBLE	long double
MPI_PACKED	
MPI_SHORT	short
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long
MPI_UNSIGNED_SHORT	unsigned short

14.2.2 一个基于 MPICH 实现多节点并行的例子

本节使用的例子详见"MPI_demo/mat_inc"目录。该程序的主要功能是给一片连续的内存中不同的位置赋不同的值。

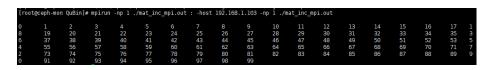
我有两台刀片服务器 node02(内网地址 192.168.1.101)、node03(内网地址 192.168.1.103),采用粗粒度的域分解方式实现任务划分,两台机器分别处理内存的前半段和后半段。两台机器的 MPI 软件版本必须相同,我用的是MPICH3。基本的 MPICH 软件包只包含 mpirun、mpiexec,用于执行 MPI 程序;编译采用 MPI 程序还需要安装 MPICH-devel 开发者软件包,包含 mpicc编译器。两台刀片服务器的操作系统皆为 CentOS(具体版本未知),互连方式为 infiniteband。

以下是步骤:

- 1. 首先确保 node02 与 node03 之间可以免密 ssh 登录。使用 "ssh-keygen -t [rsa] [target host]" 命令生成 ssh-key, 储存在 id_rsa、id_rsa.pub 文件中。
- 2. 将用 mpicc 编译器编译好的程序分别拷贝到两个节点的同一相对目录下 (mpirun 命令会携带主节点的环境参数到从节点),注意代码中需要有根据进程号进行域分解的代码。
- 3. 在主节点(我以 node02 为主节点)中通过"mpirun -np [主节点进程数] [可执行文件]: -host [从节点地址] -np [从节点进程数] [可执行文件]"执行程序。也可以通过 machinelist 添加更多的节点。

以下是通过执行命令 "mpirun -np 1./mat_inc_mpi.out: -host 192.168.1.103 -np 1./mat_inc_mpi.out" 产的生输出结果,该结果与对应串行程序的输出相

同,验证了并行程序的正确性。



14.2.3 群体通信