

计算科学与工程中的并行编程技术

Parallel Programming Technology in Computational Science and Engineering

胡志军 清华大学计算机系

email: huzh@tsinghua.edu.cn (一小时内回复)

phone: 62782530

http://hpclab.cs.tsinghua.edu.cn/~huzh

主要内容

并行编程基础 50%

综合语言与库 50%

并行编程基础 (胡志军 80%)

论文 (EISCI, 清华期刊与会议论文)

同MPI, OMP, CUDA, OpenCL 作为主要或重要通信与计算项目

参考文献

MPI - the Complete Reference

Using MPI: portable parallel programming with the MPI

Using MPI-2: advanced features of the message-passing interface

高性能并行编程技术 - MPI 并行编程技术 (胡志军)

MPI-forum.org

并行编程技术
并行编程技术

OpenMP: <http://openmp.org>

CUDA / OpenCL: <https://developer.nvidia.com/cuda-zone>

并行编程技术

Top500.org

Graph500.org

网络编程技术

www.netlib.org

(最新的、最权威的、都是开源免费的、总是会比自己写的要好)

课程目标:

并行编程技术的基本概念

并行编程技术的基本概念

并行编程技术的基本概念

并行编程技术的基本概念

并行编程技术的基本概念

最重要、也是最难学的

现实生活中的问题: 有串行的事情, 也有并行的事情

而以往的编程训练都是串行的



理想不高于现实世界的并行编程

下次课程: 应用并行编程的知识 (计算机、化学、生物...)

domain-specific optimization (2019图灵奖得主)

课程主要内容

MPI (80%) - 应用最广

其他技术 (20%) - CUDA / GPU, HPC, OpenMP, VLA...

本课程主要涉及并行编程

并行编程技术: 并行编程技术 (绝对统治地位)

多线程/多线程化、加速器等。

因为多线程编程具有并行性, 这也是并行编程技术

cluster 可视为并行编程, 机器并行编程

Lec 1 2019年9月11日
GB-208 THU

并行编程技术的发展

1986, Hans. Meuer 提出并行编程技术

Top500

1993年, 开始发布世界Top500 (也就是并行编程技术)

现在, 并行编程技术已经 LINPACK BENCHMARK 的排序

例如: Gauss 并行编程技术

目前并行编程: Summit > Sunway Taihu Light > Tianhe 2 > Titan > Sequoia > K Computer

> Tianhe 1A > Jaguar > Roadrunner > Bluegene > The earth Simulator

> ASCI White ... > CM-5 (1993年位列第一)

厂商 HPE (HP + GI) + Others 10%

应用与系统 Others 92.6%

工业界 58.4%

学术界 11.8%

Asia 51.8%

Europe 19.6%

Cyber + Ethernet 54.4%

Infinitband 24.6%

Custom 10.4%

Cluster 90.6%

MPP 9.4%

China 22%

USA 11%

日本 2%

Linux 48.6%

其他 (各种 UNIX, LINUX) 51.4%

moore定律: 每隔18个月, 集成芯片的晶体管数量上升一倍

现在遇到的挑战 (都是大问题): 功耗、散热、互连带宽、由此产生的并行技术

如何能发挥这些芯片的能力

并行编程技术

大数据与云计算

人才

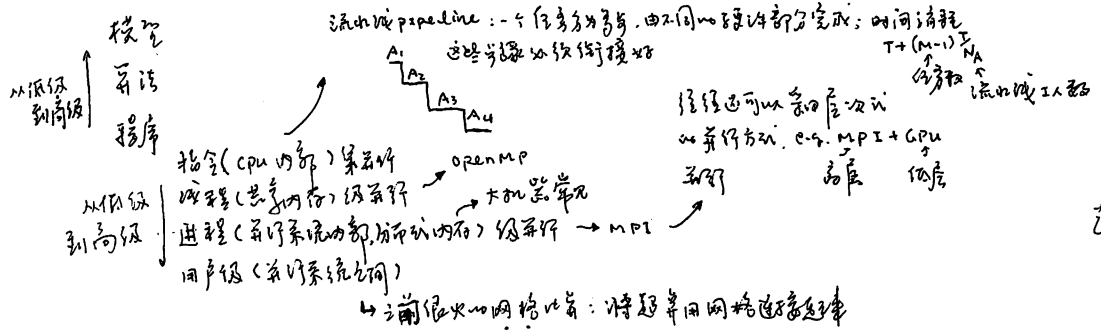
2016, 2017 Gordon Bell Award (面向应用)

Atmosphere

(Sunway)

Earthquake (Sunway)

实现并行计算的几个层次

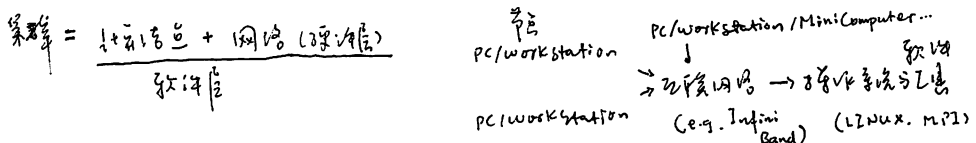


最普遍并行计算 model

消息传递 (更唯一) MPI 是消息传递的编程标准 model. 所有系统都支持
也是 TOP500 排行榜的基础

数据并行 (更高级, 更容易一点)
有一定的限制, 不可以数据并行

集群系统组成



集群系统的特点: 价格便宜, 所有组件 node/network 都是批量生产的通用部件
研发周期短, 容易实现 (设备节点及网络都是通用部件)
性能仍然在提升
体现了后摩尔定律的并行和软件

我国高性能计算面临的问题:

硬件基础薄弱 (靠进口/自主研发与生产能力)
软件投入不足 (软件研发作用, 应用, model, 算法, 设计, 开发等都要结合)
企业投入不足 (Made in China / Designed in China)
国家投入不足 高附加值产业才需要设计, 计算需求

我们的工作

集群式系统的建造 (Deep Super-21C)
163,500, 2003年, 1.5 TFlops 相当于一台显卡
256 核 Intel P4

Cluster 系统
Cluster Software System
Cluster Algorithm
TRVIA

通信中数据并行与串行 (数据并行更有优势)
问题分解是并行 OS 干预, 做网络通信

网络计算 Grid Computation
用了 Power Grid 电力网的意义

2000 年最大的技术 (类似今天的 AI)
很多超算中心, 联合起来, 协同解决更大的问题
不同单位, 管理方式完全不同, 很难集成 -> 太复杂
将一个公司内的计算资源集成 -> 可行, 即今天的 "云计算"

总结: 高性能计算是

先造基础设施, 先造软件与工程活动所需的基础
国家安全, 科学发现, 产业限制
是科学研究的第三种方法 (理论, 实验, 计算) 和科学发现的新模式
Bridge

深度学习: 给出计算方法和求解的近似解
并给出近似误差的估计

并行计算的一些概念:

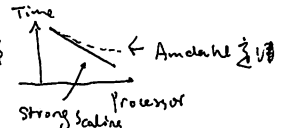
加速比 (speed-up)
Amdahl's law
若有一任务, 总执行时间为 T , 可并行部分为 α , 则 $T_1 > T_2$ 时, $Speedup = \frac{T_1}{T_2} = \frac{1}{\alpha + (1-\alpha)/n}$

即加速比 $Speedup = \frac{1}{1-\alpha + \frac{\alpha}{n}}$
当 $n \rightarrow \infty$, $Speedup = \frac{1}{1-\alpha}$

比如 $\alpha = 0.5$, 则 $Speedup = 2$
即: 若可并行部分占比不大, 并行数量与资源成正比
即: 问题的 "模型" 往往是求解的问题最重要的一步

Amdahl 定律: 问题的规模是固定的

Strong scaling: 给定问题的尺寸, 运行时间与处理器数量成反比
Weak scaling: 问题尺寸变大, 处理器数量也同步变大, 运行时间基本保持不变



即: 只要符合 Amdahl 定律, 作为 Speedup 的上限, 但实际中, 可以求解更大的问题而保持时间不变

本书主要内容：
{ 本书主要内容
{ 本书主要内容
{ 本书主要内容

并行程序是否可串行执行？

自动并行化 (编译) → 编译器, 硬件 → CPU, 自动并行化以充分挖掘并行性
手动并行化改造 → 并行编程
{ 自动并行化 (编译) → 编译器, 硬件 → CPU, 自动并行化以充分挖掘并行性
{ 手动并行化改造 → 并行编程

并行编程的挑战

并行编程的挑战：
{ 并行编程的挑战：
{ 并行编程的挑战：
{ 并行编程的挑战：

并行程序相对串行程序最大的不同

各单元如何并行执行, 如何相互协调
(embarrassingly parallel)
简单, 不需要交互与同步的并行问题
因而通信(协调)是并行单元执行是并行编程的关键问题

并行程序不一定能并行执行

(可以模拟串行, 只是效果不好而已)

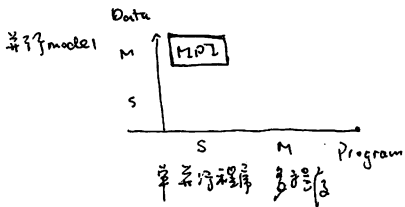
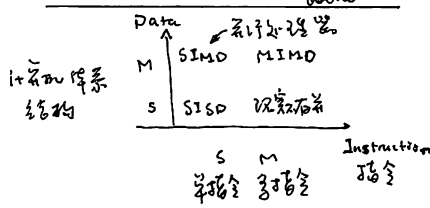
串行程序并行化是否
有意义从应用角度看

什么是 MPI (Message passing Interface) ~ 提供程序间通信交互的方法

是消息传递编程模型
是消息传递编程模型
是消息传递编程模型

MPI 在不同语境下含义不同

MPI 指的是 SPMD 模型
Single Program Multiple data



一个语言程序, 可包含多个进程, 运行时可被看作
现出不同的行为
MPI 的最高境界: 并行编程

MPI 的消息

通信的方式: 点对点通信 (点对点通信)
有时, 消息是广播的 (组通信)
自由于管理通信
分为: 消息信封: 包括源/目的地, 标识, 通信域 (三元组)
消息内容: data
数据模型: 三元组
Source: 源地址, 数据个数
Destination: 目的地, 数据个数
Tag: 标识, 通信域 (三元组)
Communication: 通信域 (三元组)
MPI-Comm-world: 通信域 (三元组)
MPI-Comm-world: 通信域 (三元组)
MPI-Comm-world: 通信域 (三元组)
MPI-Comm-world: 通信域 (三元组)

MPI 的数据类型

是一组数, 相同或不同的基本数据类型, 以连续或不连续的方式形成数组
MPI 是强数据类型
A=B 且仅当 A, B 同类型时才允许
在 MPI 中, 发送与接收的数据类型必须匹配

MPI 预定义的基本数据类型

MPI-INTEGER	↔	int
REAL		float
CHAR		char
INT		int
BYTE		
PACKED		

三种通信方式

点到点通信 (Point-to-Point)

有且只有两方 (发送方与接收方) 参与
双方必须有匹配的通信信封和通信语言
S TAG ↔ R TAG
Recv 语言 ↔ Send 语言

组通信 (Collective Communication)

组内所有进程通参为标识完成
(有且只有, 有且只有)
广播与 Recv 匹配
都是错误
通信的调用形式都相同

单边通信 (One-side Communication / Remote memory Access)

单边通信单方法或
PUT GET (可用其他通信代替?)

90 年代的 PVM, 现在已淘汰
(Virtual machine)
引入了单边通信

MPI的概述

- 是通信的对象
- 是MPI程序运行所依赖
- 是进程内具有唯一身份的对象

基本通信函数

基本发送 MPI-SEND (buf, count, datatype, dest, tag, comm);

基本接收 MPI-RECV (buf, count, datatype, source, tag, comm, status);

注意: 对于接收函数, source 可以是 MPI-ANY-SOURCE (通配符), tag 可以是 MPI-ANY-TAG (通配符)

FORTRAN和C的调用语法不同

大小无关 有关

因此C语言中最后参数 error, 而C语言通过函数返回值返回

获取身份

MPI-Comm-Rank (comm, rank);

MPI是-4SPMD模型

通过rank获取自己的身份, 从而不同身份执行不同任务, 相互协作

标识: 具有唯一-4且不同-4

同-4进程 同-4进程, 在-4 Comm #

在-4 Comm #

已集成进进程的变量 (e.g. recv from rank-1, send to rank+1)

获取组内所有大小

MPI-Comm-Size (comm, size);

可以接收并更加通用, 当任务规模不同时, 每个进程分得的任务数也不同

初始化与结束语言

MPI-INIT ();

MPI-FINALIZE ();

MPI是-C, FORTRAN库, 需要安装库并链接

但是, MPI程序的第一条语句必须是MPI-INIT, 最后一条语句必须是MPI-FINALIZE, 不可省略, 否则会出现莫名其妙的错误

(e.g. 以语言形式提出并行机制, 从而可自动完成环境设置)

若以并行语言扩展, 则是, 则介乎两者之间

程序源代码

① Hello world

```
#include "mpi.h"
MPI-INIT ();
... (e.g. printf "hello world")
MPI-FINALIZE ();
```

Comm-rank
-size
-Get-processor-name

```
include "mpi.h"
integer rc, ierr
call MPI-INIT(ierr)
... write (*,10) myid, numprocess, processor-name
call MPI-FINALIZE(rc)
```

注意: 以不同进程数, 在不同的处理器上执行

结果不同:

C5 FORTRAN 调用形式不同, 参数不同:

② A将消息发给B和B再发

```
include "mpi.h"
init ();
comm-rank (-, &myid); id = myid;
-size ();
if (myid == 0) { MPI-SEND (id, 1, MPI-INT, 1, 100, -); }
else if (myid == 1) { MPI-RECV (&id, -, 0, 100, -); printf (~id); }
finalize ();
```

分析:

进程rank	0	1	2 ...
执行流	init myid = 0 id = 0 send (id) → finalize	init myid = 1 id = 1 recv (id), 使得 id = 0 printf (id), 使得 0 被 print 出来 finalize	init myid = 2 id = 2 finalize

也体现了MPI的“局部并行”与“全局并行”的特点 (也是具有两个特点)

从-4个进程, 程序在-4个进程, 看所有-4个进程, 则是并行

消息传递有两种编程模式: 主-从, SPMD 可以写出两种模式

程序设计方法: 如并行计算
↓
设计单个任务的并行
↓
所有任务统一表示 SPMD

MPI 并行程序设计原则

- 程序结构 - 并行技术 model: 大问题分解
- 每次 Communication 之间 wait 变量, 避免通信语言阻塞.
- 降低通信频率, 甚至是以异步交换通信

作业: 在代码中, 每个进程向 rank = myrank + k 的进程发送字符串 "hello world!", 并让 rank = myrank - k 的进程接收. 这里 rank 是进程的编号, k 是给定的常数.
收到消息后, 将字符串打印出来.

常见错误: ①

```
for (inti=0; i< rank-size; ++i) {  
    if (i == myrank) { send & receive; }  
}
```

 这是串行化思想

② 接收可使用 MPI_ANY_SOURCE 没问题, 或可使用 $source = (myrank - k + rank - size) \% rank - size;$
 $dest = (myrank + k) \% rank - size;$
指定消息的源和流面 (envelop 参数)

MPI 的安装与使用

MPI CH 简介
本书为 { 关于 Linux 的 MPI CH (重点) → 真正的并行程序都是在 Linux 上运行.
关于 Windows 系统的 MPI CH

Lec 3 2019年9月25日
6B208 THU

MPI CH Argonne 国家实验室, 密西西比州大学开发
使用最广泛, 开放使用, 可移植性最好; 与 MPI 规范同步发展, MPI CH 是变通
性能可能不是最优. (recap: MPI = 模型 + 实现 + 库实现)

包一获取: google "http://www.mpi.ch.org/downloads" (可被不兼容了)

解压 tar xzf mpi.ch.tar.gz

配置与编译 \$ cd mpi.ch

\$./configure --prefix=安装目录/mpi.ch-install
\$ make
\$ make install
编译选项不同 (C, FORTRAN), 环境不同
可在此设置
e.g. CC=icc, CXX=icpc, F77=ifort
FC=ifort --with-pm=hydra

将安装路径下 /bin 加入 PATH 环境变量

doc 文档
examples 测试程序
src 源代码
lib 库
man 手册
www 官方网站

MPI 程序的编译

`mpicc (C) / mpicc (C++) / mpif77 / mpif90 -o test test.[c/c++/f/f90]`
编译器 输出文件 输入文件

MPI 程序的运行与通信机制

如构建运行程序的必要通信机制

假设1: /etc/hosts equiv 主机建立通信机制
依次给出所信机器名或 IP 地址即可.
(equiv 意思是, 地址不发生变化)
e.g. tp1.cs.tsinghua.edu.cn 机器名
tp2.cs.tsinghua.edu.cn (常见, 因为方便)
192.168.166.211

假设2: 通过 .rhosts 文件建立通信机制
(以用户为基准建立通信机制)
依次给出机器名及用户名
e.g. tp1.cs.tsinghua.edu.cn mpi

假设3 (配合假设1, 假设2)
ssh 密钥交换

第一, 生成密钥对 \$ ssh-keygen -t rsa
(在 \$HOME/.ssh/ 目录下生成 id_rsa, id_rsa.pub)
第二, 分发公钥
\$ cat \$HOME/.ssh/id_rsa.pub | ssh username@w.x.y.z 'cat >> \$HOME/.ssh/authorized-keys'

文件准备: 所有机器都要有可执行程序

同构 → 拷贝 \$ scp cp1 tp1:/home/mpi/mpi.ch/examples/basic/.
异构 → 分别编译 源 目标 (每个机器上编译)

另外, 所编的可执行文件也要拷贝过去.

运行 no cluster, 每个机器都有 local disk

另外还有用网络连接的共享 disk, 只需要数据放在共享 disk, 即不用拷贝; 每个进程的输入输出文件可以不同 (只要 rank 命名即可)

和 PVM 的差别: PVM 可以调整进程数
而 MPI 则不行, 启动时定死
(MPI-2 是加入了一些机制, 让更多的进程可以加入)

程序运行

\$ mpiexec -n N executable-file (单机运行, 调用缺省配置文件)
\$ mpiexec -f mf -nn executable-file
\$ mpirun -machinefile mf -np N executable-file (无分布)

帮助命令 \$ mpi man

配置文件如 tp1 ← 运行进程 at 机器 1
tp2:4 ← 运行 4 进程 at 机器 2
tp3:2
...

注意: 若 machine file 中指定的总的进程数少于 -n 指定的进程数。
则将以循环的方式重新遍历 machine file 中指定的主机与进程数以此创建新的进程。
注意: 在机器 Titan (曾经一起霸主), 用 aprun 而不是 mpirun

基于 OpenPBS 的执行

http://www.pbsworks.com/
开源的作业管理系统, 制定作业策略; 用户向队列提交请求
在队列上不允许登录 e.g. 作业队列中, 每次循环, 我可以满足需求 (核, 内存需求)
并结束
qsub / qdel / qstat
分别用于提交, 删除, 查询提交 信息

qsub 服务器;

```
qsub -V -l nodes=4:ppn=16  
-e 错误输出文件  
-o 标准输出文件  
-l walltime=99:00:00 时间限制, 运行时间;  
-v 有输入输出, 在节点中资源
```

```
cd PBS_O_WORKDIR  
mpirun -np 64 -machinefile $PBS_NODEFILE $PBS_O_WORKDIR/bk.exe
```

e.g. 神威 为谢尔曼 使用。
通过 easyconnect 软件, 连接 VCN 至服务器 shell。
然后, 通过 ssh 登录到另一台服务器并结束

提交作业命令: qsub -I -q gsw-share -b -m 1 -N 16 -np 4 -j ... -share-size ...
-I 交互式
-q 队列名
-b 有输出
-m 每个节点
-N 节点数
-np 进程数
-j 作业文件
-share-size 共享大小
-host-stack 1024
-csp 64 ~ /swvasp > log 2018
每个节点 64 个内核

在 Windows 下使用 MPI-GH

下载与同样的网站, 解压, 注册
mpirun.exe (注册, 与操作系统一致)

```
执行 mpirun -n 2 d:\helloworld.exe  
-hosts 2 m1 m2 d:\helloworld.exe (2 节点, 节点 m1, m2)  
-configfile d:\Config  
-n 1 -host m1 D:\helloworld.exe  
-n 1 -host m2 D:\helloworld.exe
```

也可使用 MS MPI STUDIO
Developer

练习: 提交, 设置一种 MPI, 运行一个 MPI 程序; 与老师在运行报告
机器
os.
环境, 设置.
编译
运行
所有出现的错误和解决方案

MPI的数据类型

Lec 4 2019年10月9日
GB208 THU

引入原因: 抽象
易编程
移植性好

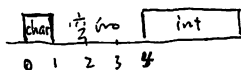
数据类型举例: 非常灵活, 可描述任何抽象概念

分类: 预定义数据类型 (MPI-INTEGER, MPI-CHAR, MPI-SHORT, ...)
自定义数据类型

自定义数据类型以预定义数据类型为基础, 最终是以基本数据类型为最终表示形式
是若干个基本数据类型在内存中连续存放的排列

对齐问题: 若 $size(T) = s$, 那么使 $address_{0_s} = 0$ 的变量, 称为对齐变量;

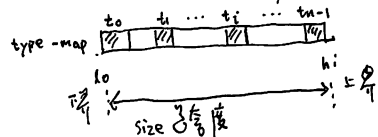
例如 $size(int) = 4$, 那么 int 只应放在地址为4的倍数。



MPI是强类型语言. send/receive 必须对数据类型严格匹配。

通过匹配是指, 若两种 datatype 最终分解为基本数据类型相同, 则可匹配。

类型图: 类型图 = { <基类型, 偏移>, <基类型, 偏移>, ... }



1. 连续类型复制:

$MPI_TYPE_CONTIGUOUS(count, oldtype, newtype);$
将 $oldtype$ 重复 $count$ 次, 作为一个新的 $newtype$ 。

例如: 有非常连续发送数据的行, 可定义为一行作为一个新的 datatype。

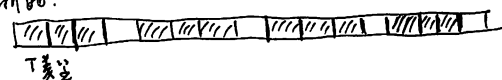
每次发送一行即可。 $MPI_TYPE_CONTIGUOUS(1000, MPI_FLOAT, 8 \times C-R);$
 $MPI_SEND(8 \times C-R, 1, C-R, right, tag, comm);$

注意: 连续重复多次, 中间也可以有空隙 (数据的 $size$ 和 $extent$ 不同时)

2. 任意类型的数据:

$MPI_TYPE_VECTOR(count, blocklength, stride, oldtype, newtype);$
类型重复 $count$ 块, 块内有 $blocklength$ 个 $oldtype$ 连续排列, 块之间有间隔 $stride$ (单位为 $oldtype$ 的长度)

例如:



$count = 4, blocklength = 3, stride = 4, oldtype = T$

注意: 也可知道, $vector$ 比 $contiguous$ 更通用。

$contiguous(C, 0, N) \Leftrightarrow vector(1, C, 0, 0, N)$
 $stride$ 为偏移量

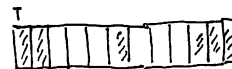
例如: 在 FORTRAN 中 $matrix(0:n-1)$ 而 C 语言中 $matrix(0:n)$ (注意: 对 $extent$ 计算全部包含的字节)

$MPI_TYPE_VECTOR(1000, 1, 1000, MPI_FLOAT, 8 \times C-COL);$

$MPI_SEND(8 \times C-COL, 1, C-COL, up, tag, comm);$

3. $MPI_TYPE_INDEXED(count, array_of_blocklengths, array_of_displacements, oldtype, newtype);$
 $oldtype$ 可被分成许多块, 每块可不等长度, 块的数量为 $count$
也可不如此 (用 $oldtype$ 的长度为单元)

例如:



$count = 3, array_of_blocklengths = \{2, 1, 3\}, array_of_displacements = \{0, 5, 9\}$
 $oldtype = T$

注意: $indexed$ 比 $vector$ 更一般

$vector(C, L, S, 0, N); \Leftrightarrow indexed(C, \{L, \dots\}, \{0, S, 2 \times S, 3 \times S, \dots, (C-1) \times S\}, 0, N);$

例如: 定义 5×5 的 $double$ (每行 5 个)

$double a[100][100];$
 $int disp[100], blocklen[100]; int i;$
 $for(i=0; i<100; ++i) \{$
 $disp[i] = 100 \times i;$
 $blocklen[i] = 100 - i; \}$

$MPI_TYPE_INDEXED(100, blocklen, disp, MPI_DOUBLE, 8 \times upper);$

$MPI_SEND(a, 1, upper, dest, tag, comm);$ 发送每行数据

4. 结构数据类型

MPI-Type-struct (Count, arr-of-blocklengths, arr-of-displacements, array-of-types, newtype);

将若干个不同数据类型组合，成为新的类型。



例如: 3, {2, 1, 3}, {0, 6, 18}
count, arr-of-blocklengths, arr-of-displ → 注意: 以字节为单位。类型是 MPI-Int,

注意: struct 是比 indexed 更一般的类型。

indexed (C, {L1, L2, ..., Lc}, {D1, D2, ..., Dc}, 0, N);

struct (C, {D1, D2, ..., Dc}, {0, 1, ..., N});

例如: struct part-t { char class; double d[6]; char b[7]; } 对应的类型

MPI-DataType type[3] = {MPI-CHAR, MPI-DOUBLE, MPI-CHAR};

int blocklength[3] = {1, 6, 7};

MPI-Int displ[3] = {0, sizeof(double), 7 * sizeof(double)};

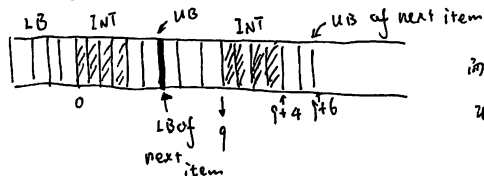
MPI-Type-struct (3, blocklength, displ, type, &newtype);

MPI-LB, MPI-UB 是4字节对齐, 左边界为0, 指定数据类型在 Lo, hi bound 内

MPI-Type-structure (3, {1, 1, 1}, {3, 0, 6}, {MPI-LB, MPI-INT, MPI-UB}, newtype);



即新的数据类型 extent 为 UB-LB=9; 注意: extent 用于计算连续偏移量, jump size.



而新 newtype 的块是 contiguous 的, 且是 block

新的类型通过 MPI-Type-commit (newtype) 和 MPI-Type-free (newtype) 提交与释放。

MPI-Address (location, address); 注意: 用于 MPI-Bottom 的地址;
例如: &a[1], &a[10][10]

例如: struct {int a; float b;} value;

MPI-Int displ;

MPI-Address (&value.a, &displ[0]);

MPI-Address (&value.b, &displ[1]);

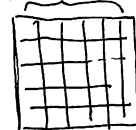
displ[1] = displ[0];

displ[0] = 0;

这样, 可以避免手动计算偏移量;

例子: 通过 datatype 来实现文本的传输 (在某些程序中, 利用 sendrecv 在间-通信传输)

Non 8014



float a[100][100], b[100][100];

MPI-Type-struct (100, 1, 100, MPI-REAL, &column);

MPI-Type-extend (MPI-REAL, &sizeext);

MPI-Type-struct (2, {1, 1}, {0, sizeext}, {column, MPI-UB}, &column);

MPI-Type-commit (&column);

MPI-Sendrecv (a, 100, column, myrank, 0, &status, b, 100 * 100, MPI-REAL, myrank, 0, &comm);

将 send 和 recv 写在一起, 直接传 100 行。
• 有 MPI 的通信, 但是是连续的, 不需要新的数据类型。
• 避免死锁

对于 1D 数组, 1D 的 type-struct, 而使用 hvector (按字节指定偏移量, 为 vector)

MPI-Type-hvector (100, 1, sizeext, column, xpose);

MPI-Type-commit (&xpose);

MPI-Sendrecv (a, 1, xpose, myrank, 0, &status, b, 100 * 100, MPI-REAL, myrank, 0, &comm);

MPI-Bottom 与结构数据类型无涉

MPI-Address (&value.a, &displ[0]);

&value.b, &displ[1];

MPI-Type-struct (2, ..., displ, {MPI-INT, MPI-FLOAT}, newtype);

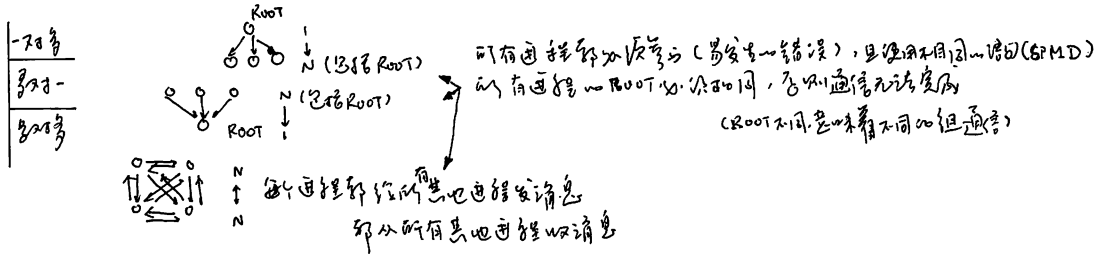
MPI-SEND (MPI-Bottom, 1, newtype, ...);

作业: 使用 MPI 的类型, 发送一文本 (M x N) 的在上三部分 (包括对角线);

MPI 组通信 { 通信
计算
同步
还会并序通信 - 一个参与优化问题, 组多进程通信
介绍
MPI-3 加入一些非阻塞的组通信, 这里先介绍

Lec 5
2019年10月16日
6B208 THU
郝志峰教授

Collective
组通信的通信功能

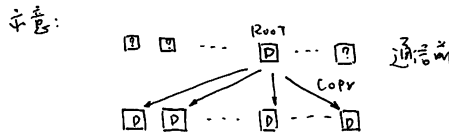


与 P-P 通信对比: P-P 通信 且 没有两个进程参与 (除 MPI-PROC-NULL 的 src / dest 时)
空进程, 通信将被忽略

- MPI-PROC-NULL 有初始值的方法. e.g. in Cartesian topology 中, 若无初始值, 边界进程值 MPI-PROC-NULL 的 src / dest, 所有进程在 MPI 调用将一致.

广播 MPI-BCAST (buffer, count, datatype, root, comm);

将 root 进程的数据广播给其他所有进程 (envelope 没有 tag, 靠调用顺序匹配)



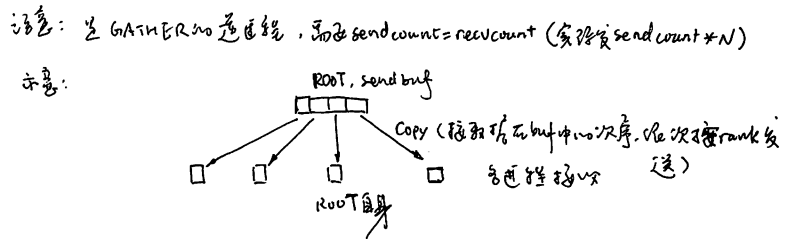
```
e.g. printf("n = %d\n", n);
if (myid == 0)
    scanf("%d", &n);
MPI_Bcast(&n, 1, MPI_INT, 0, comm); // 注意, bcast 的含义,
// ROOT 进程执行发送
printf("after bcast, n = %d\n", n); // 其他进程接收数据

o.g. test_array = (int*) malloc(sizeof(int) * nproc);
for (int i = 0; i < nproc; ++i) {
    test_array[i] = rank + i; // 进程 0: 0, 1, 2, 3, ...
    // 进程 1: 1, 2, 3, 4, ...
}
MPI_Bcast(test_array, nproc, MPI_INT, 0, MPI_COMM_WORLD);
// 所有进程: 0, 1, 2, 3, ...
```

收集 MPI-GATHER (sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm);

注意: $recvcount == sendcount$, 即所有进程发送相同的数据 type 的数据, 且个数是一样, 且接收进程接收的数据个数相同 (开始没写接收了 ROOT 的 $recvbuf$ 是多大; 非 ROOT 的 $recvcount$ 也是 $recvcount * N$);
注意: 所有进程都从 sendbuf 中 Copy 数据到 recvbuf 中 (接收区有序存放, 按在 comm 中的 rank 顺序)

散放 MPI-SCATTER (sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm);



注意: 与 BCAST 不同, 这里所有进程收到的数据不同

组收集 MPI-ALLGATHER (sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, comm);

注意: 每个进程发一个数据
所有进程都得到 gather 的结果, 所以设 root 域.
e.g. int table[MAXPROCESS][MAXPROCESS];
blocksize = MAXPROCESS / size;
sendcount = blocksize * MAXPROCESS;
recvcount = sendcount;
for (i = rank * blocksize; i < (rank + 1) * blocksize; ++i)
 for (j = 0; j < MAXPROCESS; ++j) // 计算负责的行号
 table[i][j] = rank + i;
MPI_Allgather(&table[rank * blocksize][0], sendcount, MPI_INT, table[0][0], recvcount, MPI_INT, MPI_COMM_WORLD);
gather 后的结果 table[0][0], recvcount, MPI_INT, MPI_COMM_WORLD;
注意: 0 号负责填充
所有进程获取其他进程的
数据
MAXPROCESS = 5
size = 4

问题: 接收缓冲区与发送缓冲区有重叠, 可以这样吗?

MPI-ALLGATHER的几种实现方法

环型交换算法 (RING) : 交换 $N-1$ 次

- 从最近以远程之间通信, 局部性好; 但是通信次数比较多

递归倍增 (Recursive doubling) : 总高为 $\log_2 P$ 的数据倍增树

- 子高为 2 的 k 次幂, 由 k 个依次交换数据
- 在 k 步, 子高为 2^{k-1} 的两个进程交换数据; 每一步, 数据为 2^{k-1} 个块
- 局部性好, 但通信次数较多

Bruck 算法

第 i 个进程, k 步时, 发送 2^k 个数据到 $\text{mod}((i - \exp(2, k)), P)$ 进程
从 $\text{mod}((i + \exp(2, k)), P)$ 接收数据

邻居交换算法 (Neighbor Exchange) 算法 - 左右左右算法

01 | 23 | 45 | 67 第1次
20 | 12 | 34 | 56 第2次
01 | 23 | 45 | 67 第3次

局部性好, 并且每一步只交换相邻数据

是如交换与递归倍增的折衷

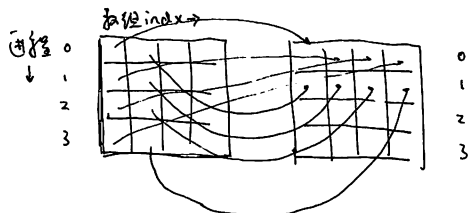
混合算法

数据通信数, 数据量

全双工 MPI-Alltoall (send buf, send count, send type, recvbuf, ^{recv count} recv type, Comm);

MPI-Alltoallv (send buf, send counts, displs, send type, recvbuf, recvcounts, rdispls, recv type, Comm);

相当于 N 次 Scatter; 相当于 N 次 Gather;



组通信的同步功能 MPI-Barrier (Comm);

组内所有进程都执行了清调用后, 才一起往下执行。

(调用程序, 数据缓冲池一致时, 使用)

eg. 在程序中, 变量 a 会不断的更新; 每轮更新时, 所有进程都打印当前值;

printf("myid = %d, the value of a = %d\n", myid, a);

MPI-Barrier (Comm);

a = ...;

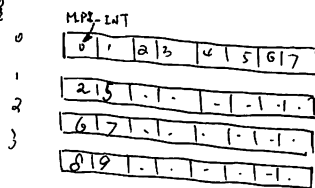
printf("myid = %d, the value of a = %d\n", myid, a);

组通信的并行功能

11.2.1 MPI-Reduce (sendbuf, recvbuf, Count, datatype, op, root, Comm);

表示 for MPI-SUM

进程



op=MPIsum

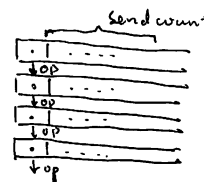
1=Root



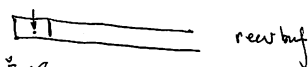
每个进程 i 个元素, 都 send 到 root 进程, 累加后放到 root 的 recv 缓冲元素。

示意图:

进程



Root



MPI数据类型: MPI_MAX, MPI_MIN, MPI_SUM, MPI_PROD, MPI_LAND, MPI_BAND, MPI_LOR, MPI_BOR, MPI_LXOR, MPI_BXOR, MPI_MAXLOC, MPI_MINLOC

Lec 6

2019年10月24日
GB208 THU
郝志彬教授

实现一：
上述语句：使用 MPI_GATHER() 以 0 进程为 ROOT。进程有 Y。然后
把 Y 最小 Y 的进程号 MIN-PROC。

然后以 MIN-PROC 广播出去，使用 BCAST()；

然后以 MIN-PROC 中的 X 广播出去，使用 BCAST()；
→ 也用 MPI_Allgather() 直接广播 Y 和 X，每个进程都自己找最小 Y 和 X。 (通信量巨大)

实现二：

用 MPI_Minloc 得到最小的 Y 的进程号，使用 MPI_Reduce()；

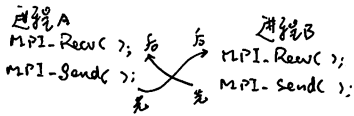
使用该进程号广播最小的 Y 对应的 X。

→ MPI_Allreduce()；

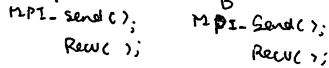
安全的 MPI 程序

能够完成所有 MPI 程序，不安全；否则是安全的

e.g. 同时收发收发，如通信锁



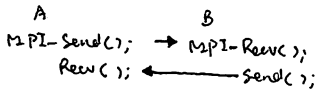
同时收发收发，可以死锁 (取决于 buffer 是否够大)



e.g. 异步收发

目前 LAN 都是 InfiniBand 或 Ethernet，旧时的集群网络都是 MyNet。
通信要求：先收线，再收发
收发要求：先收发，再收线

e.g. 收发匹配，不会死锁 (是安全的编程模式)



例如：大概都会等待，但是不会死锁

```
if (rand != 0) MPI_Send( dest = 0 );
else {
    for (i=1; i<=n-1; ++i)
        MPI_Recv( from process i );
}
```

改进的方法是使 MPI_Recv (tag=ANY_TAG, SRC=ANY_SOURCE)；

解决死锁的一种方法：MPI_SENDRECV();

MPI 的优化，且不会死锁

组调用 MPI 可以解决的问题

for() {

计算孩子；

if (B) break;

组调用 MPI；}

// 组内每个进程都执行 MPI 组调用

若有某些进程的 Break 不一致，则会导致死锁

• 对于很多串并行程序结束并行程序，有这样一个情况，e.g. B = (E < A)；

```
for ( ) {
    if (B) break;
```

• 解决死锁系统退出标志，将条件判定集中到一个进程，然后广播退出标志

```
组调用；
大家一起算；}
```

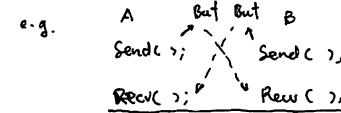
MPI 的 4 种通信模式

标准通信模式

MPI_Send();

系统决定是否缓存数据

编程者可以认为，调用结束结束，数据通信就已完成 (不一定是真的到目的进程)



开始两个 Send 都等待数据放入缓冲，然后退出调用；这样不会死锁；
若 Buffer 不够，则会死锁。

缓冲通信模式

MPI_BSEND();

使用用户提供的缓冲区 MPI_Buffer_ATTACH (buff, size);

DEATTACH (buff, size);

即如 BSEND() 时，一定放入用户缓冲区，然后退出调用；(只要 Buffer 够大，一定不致死锁)

同步通信模式

MPI_SSEND();

发送器问题同时，相应接收器问题已经解决

就绪模式

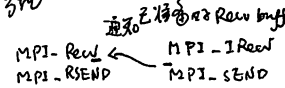
MPI_RSEND();

如数 100% 就绪，接收的 buffer 已经准备好；

• 这种中间可以取消中间通信次数，提高速度

• 但是使用起来有些危险，需要用户保证 Recv 一定是 Ready

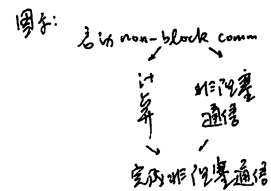
常以方法



MPI 的 非阻塞通信模式：非阻塞，发送与接收

调用 \rightarrow 返回 \rightarrow 意味着通信的完成（与之前的所有模式不同），只是留下一个发送/接收标识

作用是：计算与通信重叠（阻塞模式中，必须等待通信完成才计算）



问题来了：非阻塞 \rightarrow 发送/接收 \rightarrow 何时结束？如何知道通信是否完成？

发送 \rightarrow 非阻塞发送

```
MPI_Isend(buf, count, datatype, dest, tag, comm, request);
```

返回 \rightarrow 非阻塞发送

同时：已经开始接收，开始返回

接收：用户提供 buffer

就绪：接收已开始，才发发送

接收 \rightarrow 非阻塞接收

```
MPI_Irecv(buf, count, datatype, src, tag, comm, request);
```

测试非阻塞通信是否完成

用 request：非阻塞通信对象来描述，管理非阻塞通信

启动非阻塞通信时，创建 request

完成通信时，销毁 request

非阻塞通信的完成与测试

$MPI_Wait(request, status);$ 等待完成-移除对象 (e.g. $MPI_Isend() + MPI_Wait()$)

$MPI_Test(request, status, flag);$ 发送完成, buf 可重用 或 接收完成, 则已收到数据于 buf 中

true, 则等同于 $MPI_Wait()$ 的结果

false, 则相当无事发生

非阻塞通信 阻塞等待

等同于 $MPI_Send()$ 阻塞通信

如何使得计算与通信重叠？尽量将通信向前与

e.g. 先发送已有数据，再发送新数据

然后进行下次发送所需的数据计算

检查上次发送用发送缓冲区

e.g. 先接收旧数据 \rightarrow data 进行其他计算

检查接收，进行与接收相关计算

非阻塞通信的完成与测试

$MPI_Waitany(count, requests, index, status);$ 任何一个对象完成即返回

$MPI_Waitall(count, requests, statuses);$ 全部对象完成才返回

$MPI_Waitsome(incount, requests, outcount, array-of-indices, array-of-statuses);$ 完成 \rightarrow 对象 void

≥ 0 个对象完成即返回

完成 \rightarrow 对象 \rightarrow ids

同样也有 $Testany(count, requests, index, flag, status);$

$Testall(count, requests, flag, statuses);$ 全部完成

$Testsome(incount, requests, outcount, indices, statuses);$ ≥ 0 个对象完成即返回

例子：

```
A: MPI_Isend(Brg);
    MPI_Recv();
    MPI_Wait(Brg);

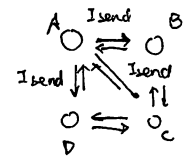
B: MPI_Isend(Brg);
    MPI_Recv();
    MPI_Wait(Brg);
```

\rightarrow 由于已经标记已发送，所以操作-立即完成

但是 \rightarrow 由于 \rightarrow send 不-是完成 \rightarrow 接收 \rightarrow 所以

还是 wait, 且 wait 是阻塞的

例子：



每个进程即的其也的进程发送消息

从其他进程接收消息

然后 wait 或 waitall

