

# Основы программирования вычислительных систем с распределенной памятью

(Занятие 3)

Игорь Николаевич Коньшин

МГУ - 03.07.2012

# Распределенная память

- Общая память → распределенная память
- Процессоры имеют свою локальную память
- Доступ к «чужой» памяти через обмены
- MPI – библиотека для обменов
  
- Большинство компьютеров из TOP-500
- Распределенная память имеется и на GPU

# Парадигма

- Доступ к нужным данным есть не всегда
- Нужны обмены
- Обмены достаточно медленные
  
- *Для общей памяти параллельность «алгоритмическая», а для распределенной – параллельность «по данным»*

# // методы линейной алгебры

- Число, вектор, матрица
- СТАТИЧЕСКОЕ распараллеливание
- Распараллеливание циклов
- Крупно-блочный параллелизм

# Оценка ускорения

$p$  - количество используемых процессоров

$T(p)$  - время решения задачи на  $p$  процессорах

$S$  - ускорение,  $S=T(1)/T(p)$

$E$  - эффективность,  $E=S/p$

$L_a$  - общее количество арифметических операций алгоритма

$L_c$  - общая длина всех обменов данными для алгоритма

$t_a$  - среднее время выполнения одной арифметической операции

$t_c$  - среднее время выполнения обмена одним числом

$T_a = L_a t_a$  - время затраченное на арифметику (вычисления)

$T_c = L_c t_c$  - время затраченное на коммуникации (обмены)

$\tau = t_c / t_a$  - общая характеристика параллельного компьютера

$L = L_a / L_c$  - общая характеристика параллельности алгоритма

$$\begin{aligned} S = S(p) &= T(1) / T(p) = T_a / (T_a/p + T_c) = p T_a / (T_a + p T_c) = p / (1 + p T_c/T_a) = \\ &= p / (1 + p (L_c t_c) / (L_a t_a)) = p / (1 + p \tau / L) \end{aligned}$$

# Идеальное распараллеливание

( $S=p$ ,  $E=1$ , linear speedup)

- (a) Сложение двух векторов:

$$Z_i = X_i + Y_i, \quad i=1, \dots, n.$$

- (b) Нормализация вектора:

$$X_i = a * X_i, \quad i=1, \dots, n, \quad a = 1 / \| X \|$$

- (c) Операция AXPY (BLAS1):

$$Z_i = a * X_i + Y_i, \quad i=1, \dots, n.$$

- *Вычисления независимы, поэтому ресурс параллелизма (максимальное возможное количество процессоров, которые можно задействовать) равен  $n$ .*

# Скалярное произведение

$$c = \text{Sum } X_i, \quad i=1, \dots, n.$$

- $L_a = n/p$ ,  $L_c \sim p-1$ ,  $L = L_a/L_c \sim n/((p-1)*p)$ ,  
 $S(L, p, \tau) = p/(1+\tau/L) = p/(1+(p-1)*p*\tau/n)$ .
- Если  $n=1000000$  и  $\tau=10$ , то  
 $S(p=1)=1$ ;  $S(p=100) \sim 90$ ;  $S(p=1000) \sim 100$ .
- `MPI_AllReduce(...)`
- *Суммирование сдвигиванием*

# MVM: плотная матрица

$$Y_i = A_{ij} * X_j, \quad i, j = 1, \dots, n$$

- (а) строчная схема хранения:

$$\begin{array}{ccc} [=] & [== == ==] & [=] \\ --- & ----- & --- \\ [=] = & [== == ==] * & [=] \\ --- & ----- & --- \\ [=] & [== == ==] & [=] \end{array}$$

- Матрица хранится по строкам, вектор точно также (и перед началом операции и после!).
- Дано,  $n$  - размерность матрицы,  $p$  - количество процессоров ( $\Pi$ ).
- Перед началом вычислений нужно собрать весь вектор на каждом  $\Pi$ .  
 $L_a = n * n / p$ ,  $L_c = (n/p) * (p-1)$ ,  $L = L_a / L_c = n / (p-1)$ ,  
 $S(L, p, \tau) = p / (1 + \tau / L) = p / (1 + (p-1) \tau / n) = S(n, p, \tau)$ .
- Если  $n=1000$  и  $\tau=10$ , то  $S(p=1)=1$ ;  $S(p=10) \sim 9$ ;  $S(p=100) \sim 50$ .

# MVM: плотная матрица

$$Y_i = A_{ij}^T * X_j, \quad i, j = 1, \dots, n$$

- (b) столбцовая схема хранения матрицы:

$$\begin{array}{ccccc} [=] & [:: & :: & ::] & [=] \\ --- & [:: & :: & ::] & --- \\ [=] = & [:: & :: & ::] & * [=] \\ --- & [:: & :: & ::] & --- \\ [=] & [:: & :: & ::] & [=] \end{array}$$

- Матрица хранится по столбцам, а вектор, естественно, по строкам.
- Для начала вычислений ничего передавать не нужно.
- Каждый процессор вычисляет частичную сумму и передает ее соответствующему П.
- Потом правда каждый П должен будет собрать частичные суммы и получить свою часть результата.
- Расчет ускорения точно такой же.

# MVM: ленточная матрица

$$\begin{array}{ccc}
 [=] & [=== & ] & [=] \\
 --- & ----- & & --- \\
 [=] & = [ & ===== & ] * [=] \\
 --- & ----- & & --- \\
 [=] & [ & === & ] & [=]
 \end{array}$$

- Пусть  $r$  - полуширина ленты.
- Перед началом вычислений нужно от двух соседних  $\Pi$  принять  $r$  компонент вектора.
- $L_a = (2r+1)*n/p$ ,  $L_c = 2r$ ,  $L = L_a/L_c = (2r+1)*n/(p*2r) \sim n/p$ ,  
 $S(L, p, \tau) = p/(1 + \tau/L) = p/(1 + p*\tau/n)$ .
- (!Note! Формула для  $S$  почти в точности та же (и почти не зависит от  $r$ ).  
 Уменьшилось количество вычислений, но уменьшилось и количество обменов.)
- Если  $n=1000000$  и  $\tau=10$ , то  $S(p=10) \sim 9.999$ ;  $S(p=1000) \sim 990$ .



# Пример MPI программы

```
#include "mpi.h"
#include "stdio.h"
int main(int argc, char** argv)
{
    int np, myid, namelen;
    char proc_name[MPI_MAX_PROCESSOR_NAME];
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &np);
    MPI_Comm_rank(MPI_COMM_WORLD, &myid);
    MPI_Get_processor_name(proc_name, &namelen);
    printf("I am %d of %d with name: %s\n",
           myid, np, proc_name);
    MPI_Finalize();
}
```

# MPI: основные функции

```
int MPI_Init( int* argc, char*** argv )
```

```
int MPI_Finalize( void )
```

```
int MPI_Comm_size(  
    MPI_Comm comm, //- идентификатор группы  
    int* size //- (OUT) размер группы  
)
```

```
int MPI_Comm_rank(  
    MPI_Comm comm, //- идентификатор группы  
    int* rank //- (OUT) номер вызывающего процесса в группе comm  
)
```

# MPI: замер времени

```
double MPI_Wtime( void )
```

```
int MPI_Barrier(  
    MPI_Comm comm //- идентификатор группы  
)
```

```
double t;  
MPI_Barrier(MPI_COMM_WORLD);  
t = MPI_Wtime();  
.....  
MPI_Barrier(MPI_COMM_WORLD);  
t = MPI_Wtime() - t;
```

# MPI\_Send

```
int MPI_Send(  
    void* buf, //- адрес начала буфера посылки сообщения  
    int count, //- число передаваемых элементов в сообщении  
    MPI_Datatype datatype, //- тип передаваемых элементов  
    int dest, //- номер процесса-получателя  
    int msgtag, //- идентификатор сообщения  
    MPI_Comm comm //- идентификатор группы  
)
```

# MPI\_Recv

```
int MPI_Recv(  
    void* buf, //- (OUT) адрес начала буфера приема сообщения  
    int count, //- максимальное число элементов в принимаемом сообщении  
    MPI_Datatype datatype, //- тип элементов принимаемого сообщения  
    int source, //- номер процесса-отправителя  
    int msgtag, //- идентификатор принимаемого сообщения  
    MPI_Comm comm, //- идентификатор группы  
    MPI_Status *status //- (OUT) параметры принятого сообщения  
)
```

# MPI\_Allreduce

```
int MPI_Allreduce(  
    void *sbuf, //- адрес начала буфера для аргументов  
    void *rbuf, //- (OUT) адрес начала буфера для результата  
    int count, //- число аргументов у каждого процесса  
    MPI_Datatype datatype, //- тип аргументов  
    MPI_Op op, //- идентификатор глобальной операции  
    MPI_Comm comm //- идентификатор группы  
)
```

```
//- op = MPI_MAX, MPI_MIN, MPI_SUM, MPI_PROD
```

# Список практических заданий (1)

Параметры // компьютера (nodes x pnp):

- максимальное число nodes
- максимальное число pnp
- максимальная память на 1 node
- максимальная память на 1 pnp

Обмены:

- внутри узла (pnp>1)
- между узлами (pnp=1)
- синхронные обмены "пинг-понг"
- синхронные обмены "чет./нечет."
- асинхронные обмены "по кольцу"
- асинхронные обмены "все всем"

# Список практических заданий (2)

Замер времени:

- время одного вызова `MPI_Wtime()`
- время одного вызова `MPI_Barrier()`
- время  $t_a$  одной арифметической операции:  $y=a*x+y$
- время  $t_c$  передачи одного числа
- коэффициент // эф-ти компьютера:  $\tau=t_c/t_a$
- измерение латентности  $t_0$  передачи данных:  $T_c=t_0+t_c*L_c$

Зависимость скорости обменов от:

- длины сообщения
- количества процессоров
- числа  $p$
- типа передачи (синхр./асинхр.)