

## 第九章 稠密矩阵运算

- ☀ 矩阵的划分
- ☀ 矩阵转置
- ☀ 矩阵-向量乘法
- ☀ 矩阵乘法

## 棋盘划分

☀  $8 \times 8$  阶矩阵,  $p=16$

P <sub>0</sub>	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
(1,0) (1,1)	(1,2) (1,3)	(1,4) (1,5)	(1,6) (1,7)
(2,0) (2,1)	(2,2) (2,3)	(2,4) (2,5)	(2,6) (2,7)
P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>	P <sub>7</sub>
(3,0) (3,1)	(3,2) (3,3)	(3,4) (3,5)	(3,6) (3,7)
(4,0) (4,1)	(4,2) (4,3)	(4,4) (4,5)	(4,6) (4,7)
P <sub>8</sub>	P <sub>9</sub>	P <sub>10</sub>	P <sub>11</sub>
(5,0) (5,1)	(5,2) (5,3)	(5,4) (5,5)	(5,6) (5,7)
(6,0) (6,1)	(6,2) (6,3)	(6,4) (6,5)	(6,6) (6,7)
P <sub>12</sub>	P <sub>13</sub>	P <sub>14</sub>	P <sub>15</sub>
(7,0) (7,1)	(7,2) (7,3)	(7,4) (7,5)	(7,6) (7,7)

P <sub>0</sub>	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
(0,0) (0,4)	(0,1) (0,5)	(0,2) (0,6)	(0,3) (0,7)
(4,0) (4,4)	(4,1) (4,5)	(4,2) (4,6)	(4,3) (4,7)
(1,0) (1,4)	(1,1) (1,5)	(1,2) (1,6)	(1,3) (1,7)
P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>	P <sub>7</sub>
(5,0) (5,4)	(5,1) (5,5)	(5,2) (5,6)	(5,3) (5,7)
(2,0) (2,4)	(2,1) (2,5)	(2,2) (2,6)	(2,3) (2,7)
P <sub>8</sub>	P <sub>9</sub>	P <sub>10</sub>	P <sub>11</sub>
(6,0) (6,4)	(6,1) (6,5)	(6,2) (6,6)	(6,3) (6,7)
(3,0) (3,4)	(3,1) (3,5)	(3,2) (3,6)	(3,3) (3,7)
P <sub>12</sub>	P <sub>13</sub>	P <sub>14</sub>	P <sub>15</sub>
(7,0) (7,4)	(7,1) (7,5)	(7,2) (7,6)	(7,3) (7,7)

(a)

(b)

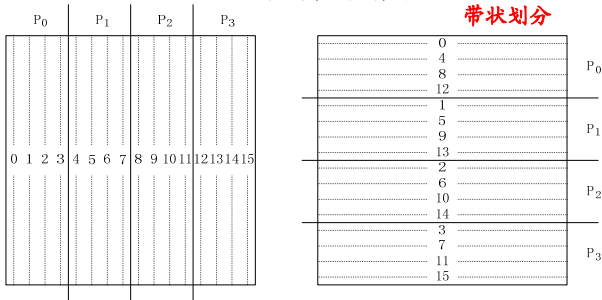
图9.2

块棋盘划分

循环棋盘划分

## 9.1 矩阵的划分

带状划分



(a)

(b)

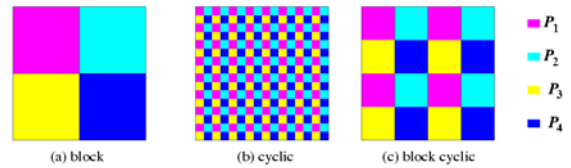
☀ 列块带状划分

图9.1

☀ 行循环带状划分

## 棋盘划分

☀ 示例:  $p=4$ ,  $16 \times 16$  矩阵的3种棋盘划分



(a) block

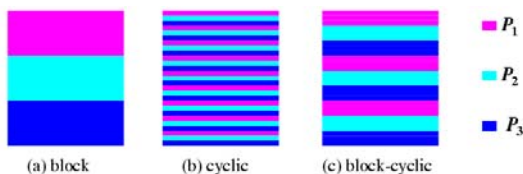
(b) cyclic

(c) block cyclic

Checkerboard mapping of a  $16 \times 16$  matrix on  $p = 2 \times 2$  processors.

## 带状划分

☀ 示例:  $p=3$ ,  $27 \times 27$  矩阵的3种带状划分



(a) block

(b) cyclic

(c) block-cyclic

Striped row-major mapping of a  $27 \times 27$  matrix on  $p = 3$  processors.

## 9.2 矩阵转置

☀ 9.2.1 棋盘划分的矩阵转置

☀ 9.2.2 带状划分的矩阵转置

## 棋盘划分的矩阵转置

### ☀ 网孔连接

◆ 情形1:  $p=n^2$ .

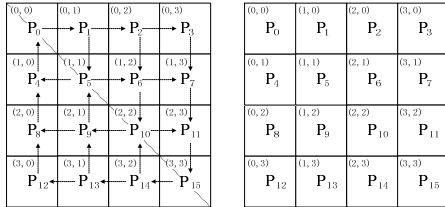
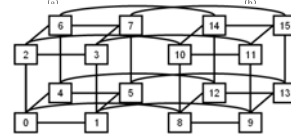
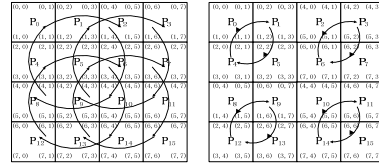


图9.3 通讯步 转置后

## 棋盘划分的矩阵转置

### ☀ 超立方连接: 示例



## 棋盘划分的矩阵转置

### ☀ 情形2: $p < n^2$ .

- 划分:  $A_{n \times n}$  划分成  $p$  个大小为  $\frac{n}{p} \times \frac{n}{p}$  子块
- 算法: ① 按 mesh 连接进行块转置 (不同处理器间) //  $2\sqrt{p}(t_s + t_w n^2/p) \dots$  通讯
- ② 进行块内转置 (同一处理器内) //  $\frac{n^2}{2p} \dots$  计算

$$T_p = \frac{n^2}{2p} + 2t_s \sqrt{p} + 2t_w n^2 / \sqrt{p} \dots \text{运行时间}$$

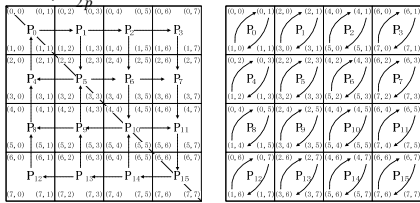


图9.4 通讯步 转置后

## 9.2.2 带状划分的矩阵转置

### ☀ 划分: $A_{n \times n}$ 分成 $p$ 个 $(n/p) \times n$ 大小的带

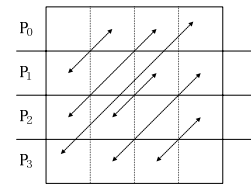


图9.7

### ☀ 算法:

- ①  $P_i$  有  $p-1$  个  $(n/p) \times (n/p)$  大小子块发送到另外  $p-1$  个处理器中;
- ② 每个处理器本地交换相应的元素

## 棋盘划分的矩阵转置

### ☀ 超立方连接

◆ 划分:  $A_{n \times n}$  划分成  $p$  个大小为  $\frac{n}{p} \times \frac{n}{p}$  子块

### ◆ 算法:

① 将  $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$  转置为  $\begin{pmatrix} A_{11} & A_{21} \\ A_{12} & A_{22} \end{pmatrix}$

② 对  $A_{ij}$  递归应用 ① 进行转置, 直至分块矩阵的元素处于同一处理器;

③ 进行同一处理器的内部转置。

### ◆ 运行时间:

$$T_p = \frac{n^2}{2p} + 2(t_s + t_w \frac{n^2}{p}) \log \sqrt{p} \quad // \text{内部转置 } \frac{n^2}{2p}, \text{选路: } 2(t_s + t_w \frac{n^2}{p}), \text{递归步: } \log \sqrt{p}$$

$$= \frac{n^2}{2p} + (t_s + t_w \frac{n^2}{p}) \log p$$

## 9.3 矩阵-向量乘法

### ☀ 9.3.1 带状划分的矩阵-向量乘法

◆  $y = A^*x$

### ☀ 9.3.2 棋盘划分的矩阵-向量乘法

## 带状划分的矩阵-向量乘法

☀划分(行带状划分):  $P_i$ 存放  $x_i$  和  $a_{i,0}, a_{i,1}, \dots, a_{i,n-1}$ , 并输出  $y_i$

☀算法: 对  $p=n$  情形

- ① 每个  $P_i$  向其他处理器播送  $x_i$  (多到多播送);
- ② 每个  $P_i$  计算;

☀注: 对  $p < n$  情形, 算法中  $P_i$  要播送  $X$  中相应的  $n/p$  个分量

(1) 超立方连接的计算时间

$$T_p = \frac{n^2}{p} + t_s \log p + \frac{n}{p} t_w (p-1) \quad // \text{前1项是乘法时间, 后2项是多到多的播送时间}$$

$$= \frac{n^2}{p} + t_s \log p + n t_w \quad // p \text{ 充分大时}$$

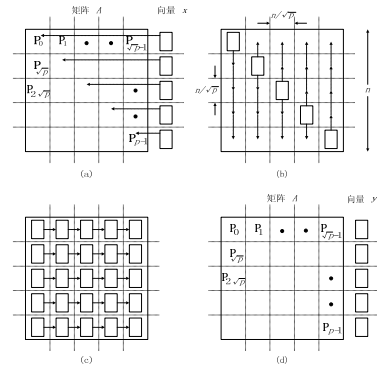
(2) 网孔连接的计算时间

$$T_p = \frac{n^2}{p} + 2(\sqrt{p}-1)t_s + \frac{n}{p} t_w (p-1)$$

$$= \frac{n^2}{p} + 2t_s(\sqrt{p}-1) + n t_w \quad // p \text{ 充分大时}$$

## 棋盘划分的矩阵-向量乘法

☀示例



## 带状划分的矩阵-向量乘法

☀示例

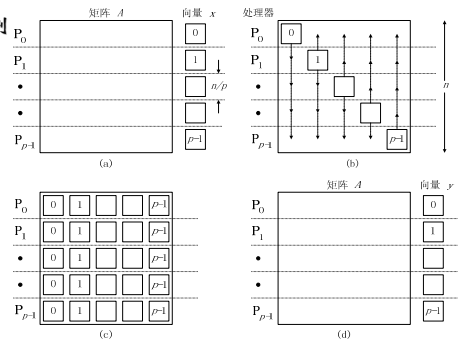


图9.8

## 带状与棋盘划分比较

☀以网孔链接为例

⊕网孔上带状划分的运行时间

$$T_p = \frac{n^2}{p} + 2t_s(\sqrt{p}-1) + n t_w \quad (9.5)$$

⊕网孔上棋盘划分的运行时间

$$T_p \approx \frac{n^2}{p} + t_s \log p + \frac{n}{\sqrt{p}} t_w \log p + 3t_s \sqrt{p} \quad (9.6)$$

☀棋盘划分要比带状划分快。

## 9.3.2 棋盘划分的矩阵-向量乘法

☀划分(块棋盘划分):  $P_{ij}$ 存放  $a_{i,j}$ ,  $x_i$ 置入  $P_{i,i}$  中

☀算法: 对  $p=n^2$  情形

- ① 每个  $P_{i,i}$  向  $P_{j,i}$  播送  $x_i$  (一到多播送);
- ② 按行方向进行乘-加与积累运算, 最后一列  $P_{i,n-1}$  收集的结果为  $y_i$ ;

☀注: 对  $p < n^2$  情形,  $p$  个处理器排成  $\sqrt{p} \times \sqrt{p}$  的二维网孔,

算法中  $P_{i,i}$  向  $P_{j,i}$  播送  $X$  中相应的  $n/\sqrt{p}$  个分量

(1) 网孔连接的计算时间  $T_p$  (CT):  $t_s + \frac{n}{\sqrt{p}} t_w + t_s \sqrt{p}$

$X$  中相应分量置入  $P_{i,i}$  的通讯时间:  $(t_s + \frac{n}{\sqrt{p}} t_w) \log \sqrt{p} + t_s(\sqrt{p}-1)$

按列一到多播送时间:  $(t_s + \frac{n}{\sqrt{p}} t_w) \log \sqrt{p} + t_s(\sqrt{p}-1)$

按行单点积累的时间:  $\therefore T_p \approx \frac{n^2}{p} + t_s \log p + \frac{n}{\sqrt{p}} t_w \log p + 3t_s \sqrt{p}$

## 9.4 矩阵乘法

☀9.4.1 简单并行分块乘法

☀9.4.2 Cannon乘法

☀9.4.3 Fox乘法

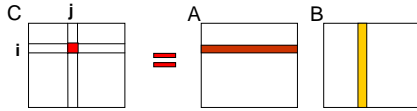
☀9.4.4 Systolic乘法

☀9.4.5 DNS乘法

## 矩阵乘法符号及定义

设  $A = (a_{ij})_{n \times n}$ ,  $B = (b_{ij})_{n \times n}$ ,  $C = (c_{ij})_{n \times n}$ ,  $C = A \times B$

$$\begin{pmatrix} c_{0,0} & c_{0,1} & \dots & c_{0,n-1} \\ c_{1,0} & c_{1,1} & \dots & c_{1,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n-1,0} & c_{n-1,1} & \dots & c_{n-1,n-1} \end{pmatrix} = \begin{pmatrix} a_{0,0} & a_{0,1} & \dots & a_{0,n-1} \\ a_{1,0} & a_{1,1} & \dots & a_{1,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n-1,0} & a_{n-1,1} & \dots & a_{n-1,n-1} \end{pmatrix} \begin{pmatrix} b_{0,0} & b_{0,1} & \dots & b_{0,n-1} \\ b_{1,0} & b_{1,1} & \dots & b_{1,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n-1,0} & b_{n-1,1} & \dots & b_{n-1,n-1} \end{pmatrix}$$



$$c_{ij} = \sum_{k=0}^{n-1} a_{ik} b_{kj}$$

A中元素的第1下标与B中元素的第2下标相一致(对准)

## 简单并行分块乘法

### 运行时间

(1)

(2) 二维环绕网孔连接:

①的时间:  $t_1 = 2(t_s + \frac{n^2}{p} t_w)(\sqrt{p}-1) = 2t_s \sqrt{p} + 2t_w \frac{n^2}{\sqrt{p}}$

②的时间  $t_2 = n^3/p$

$$\therefore T_p = \frac{n^3}{p} + 2t_s \sqrt{p} + 2t_w \frac{n^2}{\sqrt{p}}$$

### 注

- (1) 本算法的缺点是对处理器的存储要求过大  
每个处理器有  $2\sqrt{p}$  个块, 每块大小为  $n^2/p$ ,  
所以需要  $O(n^2/\sqrt{p})$ ,  $p$  个处理器共需要  $O(n^2\sqrt{p})$ ,  
是串行算法的  $\sqrt{p}$  倍

(2)  $p(n) = p$ ,  $t(n) = O(n^3/p)$ ,  $c(n) = O(n^3)$

## 矩阵乘法并行实现方法

计算结构: 二维阵列

空间对准(元素已加载到阵列中)

Cannon's, Fox's, DNS

时间对准(元素未加载到阵列中)

Systolic

$A_{0,0}$	$A_{0,1}$	$A_{0,2}$	$A_{0,3}$
$B_{0,0}$	$B_{0,1}$	$B_{0,2}$	$B_{0,3}$
$A_{1,0}$	$A_{1,1}$	$A_{1,2}$	$A_{1,3}$
$B_{1,0}$	$B_{1,1}$	$B_{1,2}$	$B_{1,3}$
$A_{2,0}$	$A_{2,1}$	$A_{2,2}$	$A_{2,3}$
$B_{2,0}$	$B_{2,1}$	$B_{2,2}$	$B_{2,3}$
$A_{3,0}$	$A_{3,1}$	$A_{3,2}$	$A_{3,3}$
$B_{3,0}$	$B_{3,1}$	$B_{3,2}$	$B_{3,3}$

## Matrix-Matrix Multiplication

Two broadcasts take time  $2(t_s \log(\sqrt{p}) + t_w(n^2/p)(\sqrt{p}-1))$

Computation requires  $\sqrt{p}$  multiplications of  $(n/\sqrt{p}) \times (n/\sqrt{p})$  submatrices

Parallel run time is approximately

$$T_p = \frac{n^3}{p} + t_s \log p + 2t_w \frac{n^2}{\sqrt{p}}$$

Algorithm is cost optimal

Isoefficiency is  $O(p^{3/2})$

— due to bandwidth term  $t_w$  and concurrency ( $p \leq n^2$  thus  $n^3 \geq p^{3/2}$ )

Major drawback of the algorithm: not memory optimal

## 简单并行分块乘法

分块: A、B和C分成  $p = \sqrt{p} \times \sqrt{p}$  的方块阵  $A_{i,j}$ ,  $B_{i,j}$ 和  $C_{i,j}$ , 大小均为  $\frac{n}{\sqrt{p}} \times \frac{n}{\sqrt{p}}$   $p$  个处理器编号为  $(P_{0,0}, \dots, P_{0,\sqrt{p}-1}, \dots, P_{\sqrt{p}-1,0}, \dots, P_{\sqrt{p}-1,\sqrt{p}-1})$ ,  $P_{i,j}$  存放  $A_{i,j}$ ,  $B_{i,j}$  和  $C_{i,j}$

算法:

① 通信: 每行处理器进行A矩阵块的多到多播送(得到  $A_{i,k}$ ,  $k=0 \sim \sqrt{p}-1$ )

每列处理器进行B矩阵块的多到多播送(得到  $B_{k,j}$ ,  $k=0 \sim \sqrt{p}-1$ )

② 乘-加运算:  $P_{i,j}$  做  $c_{ij} = \sum_{k=0}^{\sqrt{p}-1} A_{i,k} B_{k,j}$

运行时间

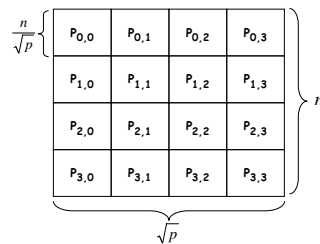
(1) 超立方连接:

① 的时间  $t_1 = 2(t_s \log \sqrt{p} + t_w \frac{n^2}{p} (\sqrt{p}-1))$

② 的时间  $t_2 = \sqrt{p} \times (\frac{n^3}{p}) = n^3/p$

## 9.4.2 Cannon乘法

分块: A、B和C分成  $p = \sqrt{p} \times \sqrt{p}$  的方块阵  $A_{i,j}$ ,  $B_{i,j}$ 和  $C_{i,j}$ , 大小均为  $\frac{n}{\sqrt{p}} \times \frac{n}{\sqrt{p}}$   $p$  个处理器编号为  $(P_{0,0}, \dots, P_{0,\sqrt{p}-1}, \dots, P_{\sqrt{p}-1,0}, \dots, P_{\sqrt{p}-1,\sqrt{p}-1})$ ,  $P_{i,j}$  存放  $A_{i,j}$ ,  $B_{i,j}$ 和  $C_{i,j}$  ( $n > p$ )



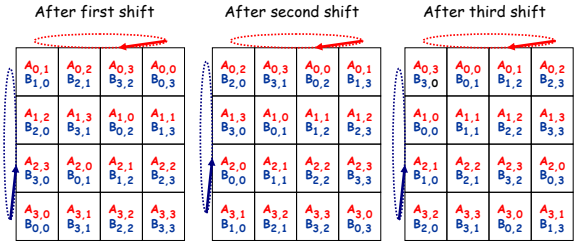
## Cannon乘法

### 算法原理 (非形式描述)

- 所有块  $A_{i,j}$  ( $0 \leq i, j \leq \sqrt{p}-1$ ) 向左循环移动  $i$  步 (按行移位);  
所有块  $B_{i,j}$  ( $0 \leq i, j \leq \sqrt{p}-1$ ) 向上循环移动  $j$  步 (按列移位);
- 所有处理器  $P_{i,j}$  做执行  $A_{i,j}$  和  $B_{i,j}$  的乘-加运算;
- A 的每个块向左循环移动一步;  
B 的每个块向上循环移动一步;
- 转②执行  $\sqrt{p}-1$  次;

## Cannon乘法

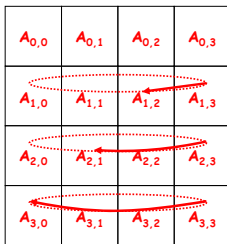
### 示例: $A_{4 \times 4}, B_{4 \times 4}, p=16$



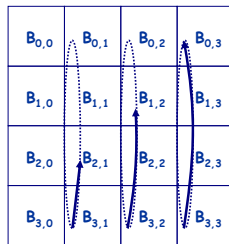
## Cannon乘法

### 示例: $A_{4 \times 4}, B_{4 \times 4}, p=16$

Initial alignment of A



Initial alignment of B



## Cannon乘法

### 算法描述: Cannon分块乘法算法

//输入:  $A_{n \times n}, B_{n \times n}$ ; 输出:  $C_{n \times n}$

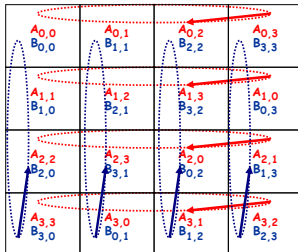
(3) for  $k=0$  to  $\sqrt{p}-1$  do  
 Begin  
 (1) for  $k=0$  to  $\sqrt{p}-1$  do  
 for all  $P_{i,j}$  par-do  
 (i) if  $i > k$  then  
 $A_{i,j} \leftarrow A_{i,(j+1) \bmod \sqrt{p}}$   
 endif  
 (ii) if  $j > k$  then  
 $B_{i,j} \leftarrow B_{(i+1) \bmod \sqrt{p}, j}$   
 endif  
 endfor  
 endfor  
 (2) for all  $P_{i,j}$  par-do  $C_{i,j} = 0$  endfor

时间分析:  
 $T_1(n) = T_1 + T_2 + T_3$   
 $= O(\sqrt{p}) + O(1) + O(\sqrt{p} \cdot (n/\sqrt{p})^3)$   
 $= O(n^3/p)$

## Cannon乘法

### 示例: $A_{4 \times 4}, B_{4 \times 4}, p=16$

A and B after initial alignment and shifts after every step



## 9.4.3 Fox乘法

### 分块: 同Cannon分块算法

### 算法原理

- $A_{i,i}$  向所在行的其他处理器进行一到多播送;
- 各处理器将收到的A块与原有的B块进行乘-加运算;
- B块向上循环移动一步;
- 如果  $A_{i,j}$  是上次第  $i$  行播送的块, 本次选择  $A_{i,(j+1) \bmod \sqrt{p}}$  向所在行的其他处理器进行一到多播送;
- 转②执行  $\sqrt{p}-1$  次;



☀ Fox (and Cannon) treatments make the following assumptions:

- ☛ The number of processes (p) is a perfect square
- ☛ The matrices to be multiplied are square of order  $n \times n$
- ☛  $\sqrt{p}$  divides  $n$  evenly

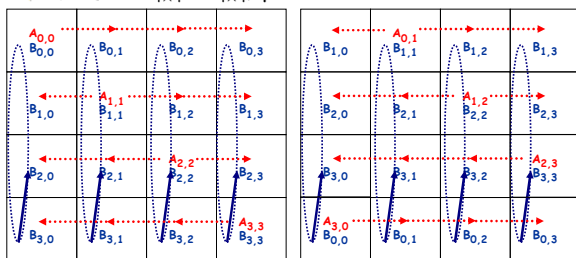
```

☀ q = sqrt(p) // number of rows, cols in processor grid
☀ // A is operand 1, B is operand 2 in A * B
☀ // C is result
☀ // i,j = process row, column
☀ // src, dest rows for rotating 'up'
☀ src = i+1 mod q;
☀ dest = i-1 mod q;
☀ for (stage = 0; stage < q; stage++) {
☛ k_bar = (i+stage) mod q;
☛ broadcast(A[i,k_bar]) to row i;
☛ C[i,j] = C[i,j] + A[i,k_bar]*B[k_bar,j]
☛ sendrecv(B[k_bar,j],src,dest);
☀ }

```

### Fox乘法

☀ 示例:  $A_{4 \times 4}$ ,  $B_{4 \times 4}$ ,  $p=16$



(a)

(b)

### 9.4 矩阵乘法

#### 9.4.1 简单并行分块乘法

#### 9.4.2 Cannon乘法

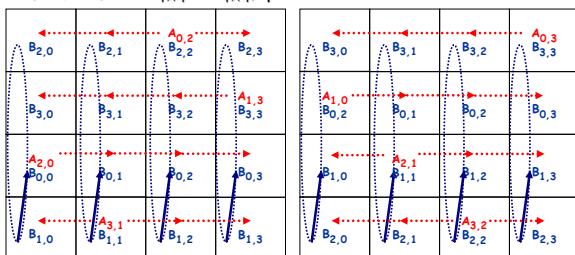
#### 9.4.3 Fox乘法

#### 9.4.4 Systolic乘法

#### 9.4.5 DNS乘法

### Fox乘法

☀ 示例:  $A_{4 \times 4}$ ,  $B_{4 \times 4}$ ,  $p=16$



(c)

(d)

### Pros and Cons of Cannon

- ☀ Local computation one call to (optimized) matrix-multiply
- ☀ Hard to generalize for
  - ☛ p not a perfect square
  - ☛ A and B not square
  - ☛ Dimensions of A, B not perfectly divisible by  $s=\sqrt{p}$
  - ☛ A and B not "aligned" in the way they are stored on processors
  - ☛ block-cyclic layouts
- ☀ Memory hog (extra copies of local matrices)

## SUMMA Algorithm

- ☀ SUMMA = Scalable Universal Matrix Multiply
- ☀ Slightly less efficient, but simpler and easier to generalize
- ☀ Presentation from van de Geijn and Watts
  - ✦ [www.netlib.org/lapack/lawns/lawn96.ps](http://www.netlib.org/lapack/lawns/lawn96.ps)
  - ✦ Similar ideas appeared many times
- ☀ Used in practice in PBLAS = Parallel BLAS
  - ✦ Basic Linear Algebra Subprograms
  - ✦ [www.netlib.org/lapack/lawns/lawn100.ps](http://www.netlib.org/lapack/lawns/lawn100.ps)

## SUMMA performance

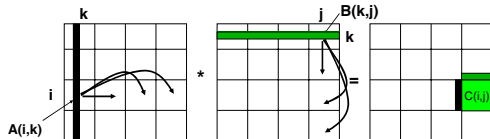
◦ To simplify analysis only, assume  $s = \sqrt{p}$

```

For k=0 to n/b-1
  for all i = 1 to s ... s = sqrt(p)
    owner of A(i,k) broadcasts it to whole processor row
    ... time = log s * (alpha + beta * b * n/s), using a tree
  for all j = 1 to s
    owner of B(k,j) broadcasts it to whole processor column
    ... time = log s * (alpha + beta * b * n/s), using a tree
  Receive A(i,k) into Acol
  Receive B(k,j) into Brow
  C_myproc = C_myproc + Acol * Brow
  ... time = 2 * (n/s) * 2 * b
  
```

◦ Total time =  $2 * n^3 / p + \alpha * \log p * n / b + \beta * \log p * n^2 / s$

## SUMMA

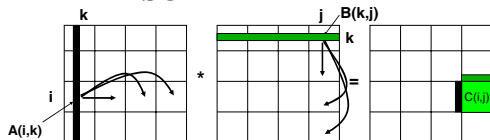


- $i, j$  represent all rows, columns owned by a processor
- $k$  is a single row or column
  - or a block of  $b$  rows or columns
- $C(i,j) = C(i,j) + \sum_k A(i,k) * B(k,j)$
- Assume a  $p_r$  by  $p_c$  processor grid ( $p_r = p_c = 4$  above)
  - Need not be square

## SUMMA performance

- Total time =  $2 * n^3 / p + \alpha * \log p * n / b + \beta * \log p * n^2 / s$
- Parallel Efficiency =  $1 / (1 + \alpha * \log p * p / (2 * b * n^2) + \beta * \log p * s / (2 * n))$
- ~Same  $\beta$  term as Cannon, except for  $\log p$  factor
  - $\log p$  grows slowly so this is ok
- Latency ( $\alpha$ ) term can be larger, depending on  $b$ 
  - When  $b=1$ , get  $\alpha * \log p * n$
  - As  $b$  grows to  $n/s$ , term shrinks to  $\alpha * \log p * s$  ( $\log p$  times Cannon)
- Temporary storage grows like  $2 * b * n / s$
- Can change  $b$  to tradeoff latency cost with memory

## SUMMA

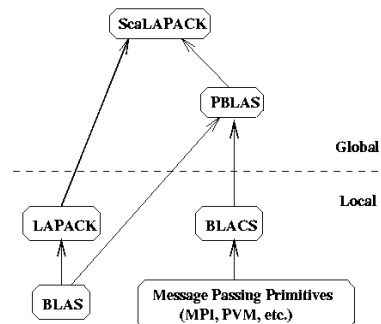


```

For k=0 to n-1 ... or n/b-1 where b is the block size
  ... = # cols in A(i,k) and # rows in B(k,j)
  for all i = 1 to p_r ... in parallel
    owner of A(i,k) broadcasts it to whole processor row
  for all j = 1 to p_c ... in parallel
    owner of B(k,j) broadcasts it to whole processor column
  Receive A(i,k) into Acol
  Receive B(k,j) into Brow
  C_myproc = C_myproc + Acol * Brow
  
```

## ScaLAPACK Parallel Library

ScaLAPACK SOFTWARE HIERARCHY



**Performance of PBLAS**

Machine	Procs	Block Size	Speed in Mflops of PDGEMM		
			2000	4000	10000
Cray T3E	4	4=2x2	1055	3070	0
		16=4x4	3830	4005	4292
		64=8x8	13456	14287	16755
IBM SP2	4	50	755	0	0
		16	2514	2850	0
		64	6205	6709	10774
Intel XP/S MP Paragon	4	32	330	0	0
		16	1233	1281	0
		64	4496	4864	5257
Berkeley NOW	4	32	463	470	0
		64	2490	2822	3450
		128	4130	5457	6647

PDGEMM = PBLAS routine for matrix multiply

Observations:  
For fixed N, as P increases Mflops increases, but less than 100% efficiency  
For fixed P, as N increases, Mflops (efficiency) rises

DGEMM = BLAS routine for matrix multiply

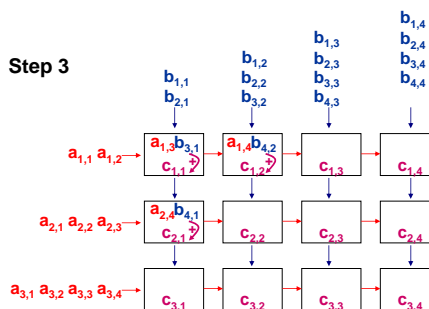
Maximum speed for PDGEMM = # Procs \* speed of DGEMM

Observations (same as above):  
Efficiency always at least 48%  
For fixed N, as P increases, efficiency drops  
For fixed P, as N increases, efficiency increases

Machine	Peak/proc	DGEMM Mflops	Efficiency = Mflops(PDGEMM)/(Procs*Mflops(DGEMM))		
			2000	4000	10000
Cray T3E	600	360	4	.73	.74
			16	.63	.70
			64	.58	.62
IBM SP2	266	200	4	.94	.89
			16	.79	.89
			64	.68	.68
Intel XP/S MP Paragon	100	90	4	.92	.89
			16	.86	.89
			64	.78	.84
Berkeley NOW	334	129	4	.90	.91
			16	.80	.68
			64	.50	.66

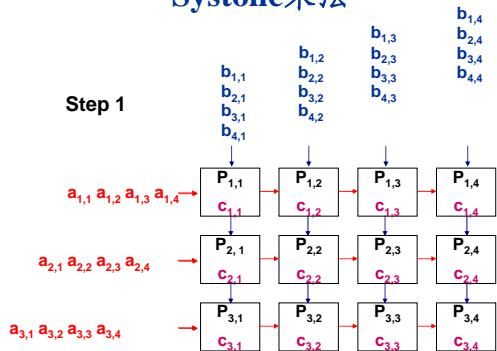
**Systolic 乘法**

Step 3



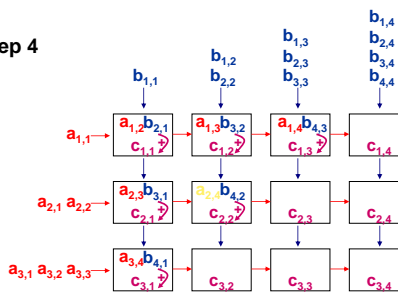
**Systolic 乘法**

Step 1



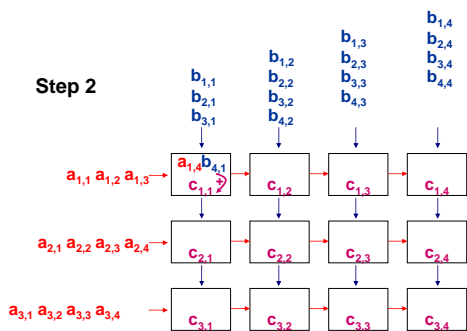
**Systolic 乘法**

Step 4



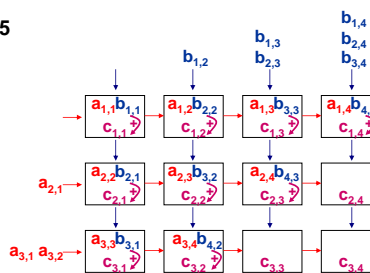
**Systolic 乘法**

Step 2



**Systolic 乘法**

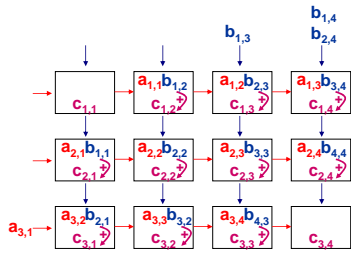
Step 5





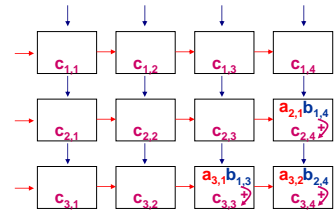
### Systolic乘法

Step 6



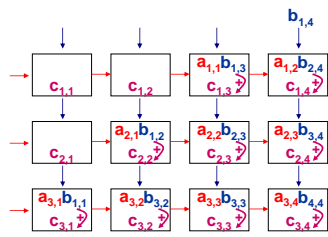
### Systolic乘法

Step 9



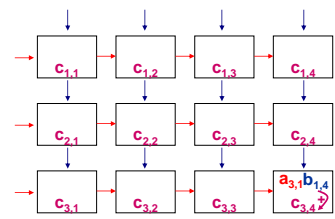
### Systolic乘法

Step 7



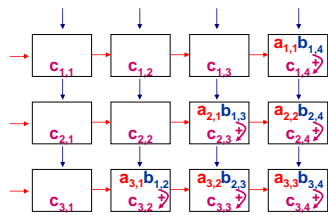
### Systolic乘法

Step 10



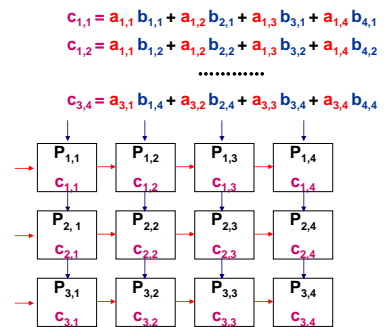
### Systolic乘法

Step 8



### Systolic乘法

Over



## Systolic乘法

```

☀ Systolic算法
//输入:  $A_{m \times n}$ ,  $B_{n \times k}$ ; 输出:  $C_{m \times k}$ 
Begin
  for i=1 to m par-do
    for j=1 to k par-do
      (i)  $c_{i,j} = 0$ 
      (ii) while  $P_{i,j}$  收到a和b时 do
         $c_{i,j} = c_{i,j} + ab$ 
        if  $i < m$  then 发送b给 $P_{i+1,j}$  endif
        if  $j < k$  then 发送a给 $P_{i,j+1}$  endif
      endwhile
    endfor
  endfor
End
    
```

## 9.4 矩阵乘法

### 9.4.1 简单并行分块乘法

### 9.4.2 Cannon乘法

### 9.4.3 Fox乘法

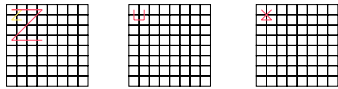
### 9.4.4 Systolic乘法

### 9.4.5 DNS乘法

## Recursive Layouts

☀ For both cache hierarchies and parallelism, recursive layouts may be useful

☀ Z-Morton, U-Morton, and X-Morton Layout



☀ Also Hilbert layout and others

☀ What about the user's view?

- ✦ Fortunately, many problems can be solved on a permutation
- ✦ Never need to actually change the user's layout

## DNS乘法

☀ 背景: 由Dekel, Nassimi和Sahni提出的SIMD-CC上的矩阵乘法, 处理器数目为 $n^2$ , 运行时间为 $O(\log n)$ , 是一种速度很快的算法。

☀ 基本思想: 通过一列和一列多的播送办法, 使得处理器 $(k,i,j)$ 拥有 $a_{i,k}$ 和 $b_{k,j}$ , 进行本地相乘, 再沿k方向进行单点积累求和, 结果存储在处理器 $(0,i,j)$ 中。

☀ 处理器编号: 处理器数 $p = n^2 = (2^2)^2 = 2^{2^2}$ , 处理器 $P_r$ 位于位置 $(k,i,j)$ ,

这里 $r = kn^2 + in + j$ , ( $0 \leq i, j, k \leq n-1$ )。位于 $(k,i,j)$ 的处理器 $P_r$ 的三个寄存器 $A_r, B_r, C_r$ 分别表示为 $A[k,i,j], B[k,i,j]$ 和 $C[k,i,j]$ , 初始时均为0。

☀ 算法: 初始时 $a_{i,j}$ 和 $b_{i,j}$ 存储于寄存器 $A[0,i,j]$ 和 $B[0,i,j]$ ;

① 数据复制: A,B同时在k维复制(一列一播送);  
A在j维复制(一列多播送); B在i维复制(一列多播送);

② 相乘运算: 所有处理器的A、B寄存器两两相乘;

③ 求和运算: 沿k方向进行单点积累求和;

## Summary of Parallel Matrix Multiplication

### ☀ 1D Layout

- ✦ Bus without broadcast - slower than serial
- ✦ Nearest neighbor communication on a ring (or bus with broadcast): Efficiency =  $1/(1 + O(p/n))$

### ☀ 2D Layout

#### ✦ Cannon

- Efficiency =  $1/(1 + O(a * (\sqrt{p}/n)^3 + b * \sqrt{p}/n))$
- Hard to generalize for general p, n, block cyclic, alignment

#### ✦ SUMMA

- Efficiency =  $1/(1 + O(a * \log p * p / (b * n^2) + b * \log p * \sqrt{p}/n))$
- Very General
- b small => less memory, lower efficiency
- b large => more memory, high efficiency

#### ✦ Recursive layouts

- ✦ Current area of research

☀ 示例

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

$$B = \begin{pmatrix} -5 & -6 \\ 7 & 8 \end{pmatrix}$$

$$\text{求 } C = A \times B$$

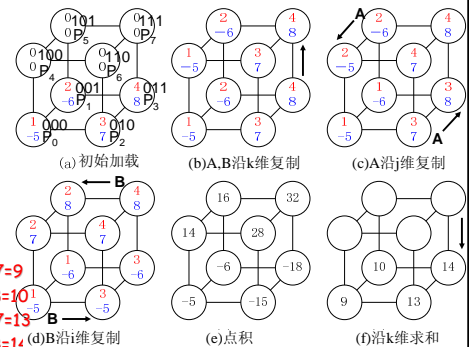


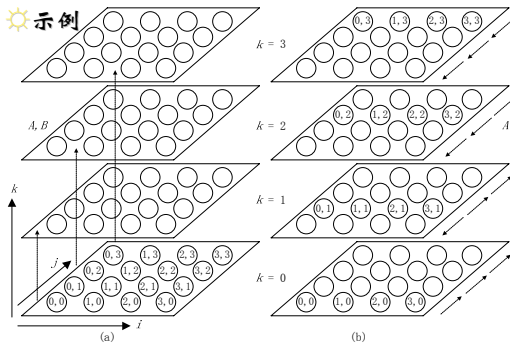
图9.12

### DNS乘法:算法描述

```

//令 $r^{(m)}$ 表示 $r$ 的第 $m$ 位取反;
// $\{p, r_m=d\}$ 表示 $r(0 \leq r \leq p-1)$ 的集合,
//这里 $r$ 的二进制第 $m$ 位为 $d$ ;
//输入:  $A_n \times n, B_n \times n$ ; 输出:  $C_n \times n$ 
begin //以 $n=2, p=8=2^3$ 为例,  $q=1, r=(r_2, r_1, r_0)$ ;
(1)for  $m=3q-1$  to  $2q$  do //按 $k$ 维复制 $A, B, m=2$ 
    for all  $r$  in  $\{p, r_m=0\}$  par-do // $r_2=0$ 的 $r$ 
        (1.1)  $A_{r^{(m)}} \leftarrow A_r$  // $A(100) \leftarrow A(000)$ 等
        (1.2)  $B_{r^{(m)}} \leftarrow B_r$  // $B(100) \leftarrow B(000)$ 等
    endfor
(2)for  $m=q-1$  to  $0$  do //按 $j$ 维复制 $A, m=0$ 
    for all  $r$  in  $\{p, r_m=r_{2q+m}\}$  par-do // $r_1=r_2$ 的 $r$ 
         $A_{r^{(m)}} \leftarrow A_r$  // $A(001) \leftarrow A(000), A(100) \leftarrow A(101)$ 
        // $A(011) \leftarrow A(010), A(110) \leftarrow A(111)$ 
    endfor
(3)for  $m=2q-1$  to  $q$  do //按 $i$ 维复制 $B, m=1$ 
    for all  $r$  in  $\{p, r_m=r_{q+m}\}$  par-do
        // $r_2=r_1$ 的 $r$ 
         $B_{r^{(m)}} \leftarrow B_r$  // $B(010) \leftarrow B(000), B(100) \leftarrow B(110)$ 
        // $B(011) \leftarrow B(001), B(101) \leftarrow B(111)$ 
    endfor
(4)for  $r=0$  to  $p-1$  par-do //相乘, all  $P_r$ 
     $C_r = A_r \times B_r$ 
(5)for  $m=2q$  to  $3q-1$  do //求和,  $m=2$ 
    for  $r=0$  to  $p-1$  par-do
         $C_r = C_r + C_{r^{(m)}}$ 
    endfor
endfor
end
    
```

### DNS乘法



### DNS乘法

