



## Matrix Algorithms: Introduction Due to their regular structure, parallel computations involving matrices and vectors readily lend themselves to data-decomposition. Typical algorithms rely on input, output, or intermediate data decomposition. Most algorithms use one- and two-dimensional block, cyclic, and block-cyclic partitions for parallel processing. The run-time performance of such algorithms depends on

- The run-time performance of such algorithms depends on the amount of overheads incurred as compared to the computation workload.
- As a rule of thumb, good speedup can be achieved if the computation granularity is able to outweigh the overheads such as the communication cost, consolidation cost algorithm penalty, data packaging, etc.































## Matrix-Matrix Multiplication: Cannon's Algorithm

- Align the blocks of *A* and *B* in such a way that each process multiplies its local submatrices. This is done by shifting all submatrices *A<sub>ij</sub>* to the left (with wraparound) by *i* steps and all submatrices *B<sub>ij</sub>* up (with wraparound) by *j* steps.
- Perform local block multiplication.
- Each block of *A* moves one step left and each block of *B* moves one step up (again with wraparound).
- Perform next block multiplication, add to partial result, repeat until all √p blocks have been multiplied.



B <sub>0,0</sub>	B <sub>0,1</sub>	B <sub>0,2</sub>	B <sub>0,3</sub>	
B <sub>1,0</sub>	B <sub>1,1</sub>	B <sub>1,3</sub>	B <sub>1,3</sub>	
B <sub>2,0</sub>	B <sub>2,1</sub>	B <sub>2,2</sub>	B <sub>2,3</sub>	
B <sub>3,0</sub>	в <sub>3,1</sub>	B <sub>3,2</sub>	B <sub>3,3</sub>	

(b) Initial alignment of B





·			2
	4		•
A.2	A <sub>0.3</sub>	-Aaa	A <sub>0.1</sub>
B <sub>2,0</sub>	B <sub>3,1</sub>	B <sub>0,2</sub>	B <sub>1,3</sub>
A1.3	A <sub>1.0</sub>	-A <sub>1.1</sub> -	A <sub>1.2</sub>
B <sub>3,0</sub>	B <sub>0,1</sub>	B <sub>1,2</sub>	B <sub>2,3</sub>
A2.0	A <sub>2.1</sub>	-A22	A2.3
B <sub>0,0</sub>	B <sub>1,1</sub>	B <sub>2,2</sub>	B <sub>3,3</sub>
A.,	A <sub>3.2</sub>	+A <sub>13</sub> -	A
B <sub>1,0</sub>	B <sub>2,1</sub>	B <sub>3,2</sub>	B <sub>0,3</sub>
	( -		(

(e) Submatrix locations after second shift (f) Submatrix locations after third shift

-1,0	-2,1	3,2	-0,5					
A <sub>1,2</sub>	-A <sub>1,3</sub>	A <sub>1,0</sub>	A <sub>1,1</sub>					
B <sub>2,0</sub>	B <sub>3,1</sub>	B <sub>0,2</sub>	B <sub>1,3</sub>					
A <sub>2,3</sub>	-A <sub>2,0</sub>	A <sub>2,1</sub>	+A <sub>2,2</sub>					
B <sub>3,0</sub>	B <sub>0,1</sub>	B <sub>1,2</sub>	B <sub>2,3</sub>					
A <sub>3,0</sub>	-A <sub>3,1</sub>	A <sub>3,2</sub>	-A <sub>3,3</sub>					
B <sub>0,0</sub>	B <sub>1,1</sub>	B <sub>2,2</sub>	B <sub>3,3</sub>					
Submatrix locations after first shift								
A <sub>0,3</sub> B <sub>3,0</sub>	A <sub>0,0</sub> B <sub>0,1</sub>	$\substack{A_{0,1} \\ B_{1,2}}$	A <sub>0,2</sub> B <sub>2,3</sub>					
A <sub>1,0</sub>	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>1,3</sub>					
B <sub>0,0</sub>	B <sub>1,1</sub>	B <sub>2,2</sub>	B <sub>3,3</sub>					
A <sub>2,1</sub>	A <sub>2,2</sub>	A <sub>2,3</sub>	A <sub>2,0</sub>					
B <sub>1,0</sub>	B <sub>2,1</sub>	B <sub>3,2</sub>	B <sub>0,3</sub>					
A <sub>3,2</sub>	A <sub>3,3</sub>	A <sub>3,0</sub>	A <sub>3,1</sub>					

 $(n/\sqrt{p}) \times (n/\sqrt{p})$ **Matrix-Matrix Multiplication: Cannon's Algorithm** A., In the alignment step, since the maximum distance over which a block shifts is  $\sqrt{p} - 1$ , the two shift operations require a total of  $2(t_s + t_w n^2/p)$  time. Each of the  $\sqrt{p}$  single-step shifts in the compute-and-shift phase of the algorithm takes  $t_s + t_w n^2/p$  time. The computation time for multiplying  $\sqrt{p}$  matrices of size  $(n/\sqrt{p}) \times (n/\sqrt{p})$ is  $n^{3}/p$ . (i.e.,  $\sqrt{p} \times (n/\sqrt{p})^{3}$ ) A<sub>0,0</sub> B<sub>0,1</sub> A<sub>0,1</sub> B<sub>1,2</sub> A<sub>0,2</sub> B<sub>2,3</sub> A<sub>1,0</sub> B<sub>0,0</sub> A<sub>1,1</sub> B<sub>1,1</sub> A<sub>1,2</sub> B<sub>2,2</sub> A<sub>1,3</sub> B<sub>3,3</sub> • The parallel time is approximately:  $T_P = \frac{n^3}{n} + 2\sqrt{p}t_s + 2t_w \frac{n^2}{\sqrt{p}}.$ A<sub>2,1</sub> B<sub>1,0</sub> A<sub>2,2</sub> B<sub>2,1</sub> A<sub>2,3</sub> B<sub>3,2</sub> A<sub>2,0</sub> B<sub>0,3</sub> A<sub>3,1</sub> B<sub>1,3</sub>

## 10















process P<sub>i,j,0</sub> accumulates the results of the multiplication from processes P<sub>i,j,1</sub>, ..., P<sub>i,j,n-1</sub>.

The DNS algorithm has three main communication steps: (1) moving the columns of A and the rows of B to their respective planes, (2) performing one-to-all broadcast along the j axis for A, and along the i axis for B, and (3) all-to-one reduction along the k axis. All these operations are performed within groups of n processes and take time O(log n). Thus, the parallel run time for 27 multiplying two n x n matrices using the DNS algorithm on n<sup>3</sup> processes is O(log n).



