

# Matrix-Vector Multiplication

**Problem** • multiply dense  $n \times n$  matrix  $A$  with  $n \times 1$  vector  $x$  ( $A \times x$ )

↑  
useless entries가 들어있는 상태

**output** •  $n \times 1$  result vector  $y$



**Serial complexity**

• serial run time =  $n^2 = \sim n \times n$   
rows      n operations in each row

```
1. procedure MAT_VECT ( A, x, y)
2. begin
3.   for i := 0 to n - 1 do
4.     begin
5.       y[i] := 0;
6.       for j := 0 to n - 1 do
7.         y[i] := y[i] + A[i, j] * x[j];
8.       endfor;
9.   end MAT_VECT
```

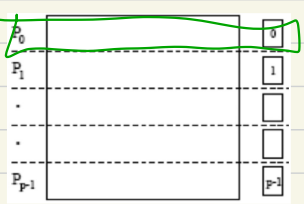
**Parallel approach**

- ✓ • row-wise 1D
- ✓ • column-wise 1D
- ✓ • 2D partitioning

# Matrix-Vector Multiplication : 1D Partitioning

Case of  $P=n$

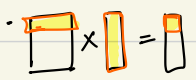
Partitioning



- each processor stores
  - row of matrix ( $n/P$ )
  - 1 element of vector

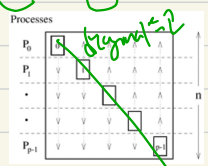
Computation

① get element  $y[i]$  from  $x$  vector and element  $a$  from matrix.

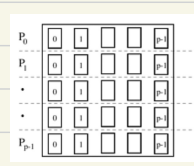


→ append to  $y$

all-to-all broadcast → each node gets entire elements of vector



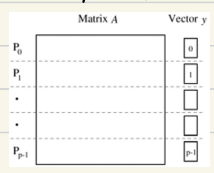
all-to-all broadcast



after broadcast

Broadcast  
Time =  $\frac{O(P^2)}{O(n)}$

② process  $P_i$  computes  $y[i] = \sum_{j=0}^{P-1} A[i,j] * x[j]$



serial work  $n^2$

Parallel Run Time

- all-to-all broadcast =  $O(n)$  time
- Computation of  $y[i] = O(n)$  time

total  $O(n^2)$  time per node

• total parallel work =  $O(P^2) = O(n^2)$

serial work의 비례  $n^2$ 에 비례함

Serial이랑 같은 speed up 일까요?

Case of pen

Pen

NLP

Com

The diagram illustrates the flow of data from a 'Case of pen' to a 'Pen' and then to 'NLP' and 'Com'. A green arrow points from 'Case of pen' to 'Pen'. Another green arrow points from 'Pen' to a green circle labeled 'NLP'. A third green arrow points from 'NLP' to a green rectangle labeled 'Com'.

- n/p complete rows of matrix A
  - n/p elements of vector x
- } per node

① 각 row own  $\vec{w}$ 를 곱하는 n/p vector elements all-to-all broadcast  
② 각 row n/p  $\times$  vec length vec row-이 곱하는 local dot product

- 2D mesh protocol

→ Computertime:  $N^2/p = (N/p) * N$

-Phase 1 :  $(t_s + t_w(n/p))(\sqrt{p}-1)$

- Phase 2:  $(t_s + t_w(n\sqrt{p}/p))(\sqrt{p}-1)$

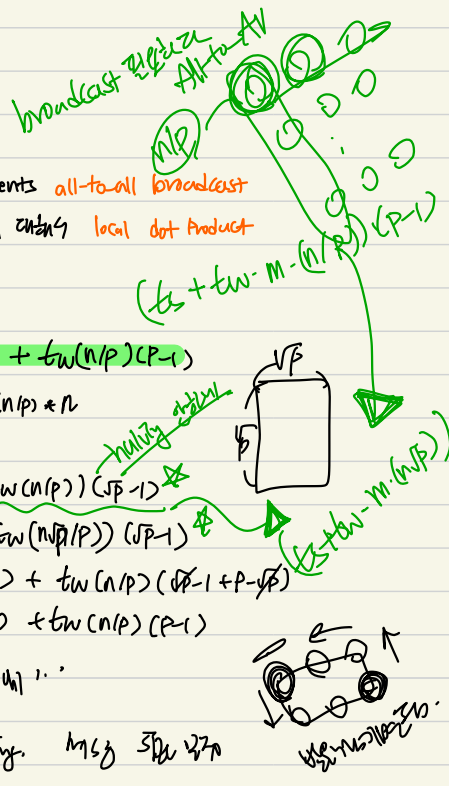
$$-total: 2ts(\sqrt{p}-1) + tw(n/p)(\sqrt{p}-1+p-\sqrt{p})$$

$$= L_f(\sqrt{p-1}) + t_w(n/p)(p-1)$$

- Ring (not Cu)

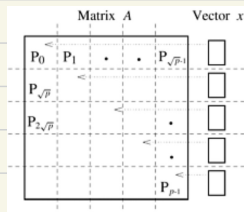
broadcast ~~renewed~~ doubling. m/s 3/27/27

$$(t_b + t_w \cdot m \cdot n / \sqrt{p}) \quad (21)$$



# Matrix-Vector Multiplication: 2D Partitioning

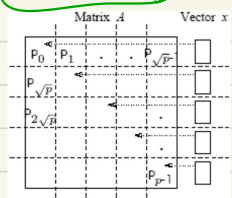
Case of  $P=n^2$  Partitioning



- each node stores
  - 1 element of matrix A
- only last node of each row holds 1 element of x

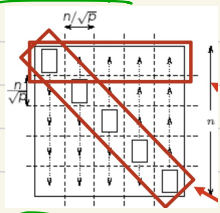
Computation

① last column node  $\leftarrow$  vector-x element  $\leftarrow$  diagonal locate  $\leftarrow$  broadcast



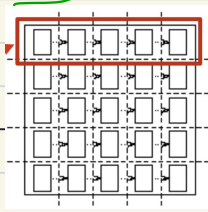
$\rightarrow$  one-to-one communication  
 $\rightarrow$  cost:  $\text{constant } O(1) \times \frac{1}{\sqrt{p}} + \frac{1}{\sqrt{p}} \times \sqrt{p} = O(1)$

② one-to-all broadcast from diagonal



•  $\frac{1}{\sqrt{p}}$  diagonal  $\times$  rowwise broadcast  $\rightarrow$  cost  
 • cost:  $O(\log n)$  time  
 recursive doubling  $\frac{1}{\sqrt{p}} \times \log n$

③ All-to-all reduction  $\rightarrow$  reduction element  $\rightarrow$  in each row to compute y[i]



$\rightarrow$  reduction element  $\rightarrow$  in each row to compute y[i]  
 • cost:  $O(\log n)$  time

Total Parallel Work  $\leftarrow$  cost  $\times$  parallel time  
 $\leftarrow O(n^2 \log n)$  ... scalable  $O(n^2)$  work, overhead is small  
 speedup of  $\frac{1}{\log n}$  is small! (ideal)  
 nodes  $\leftarrow$  (4?)

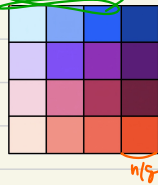
# Matrix- Matrix Multiplication

## Matrix- Matrix Multiplication

- Compute  $C = A \times B$  for  $A, B, C$   $N \times N$  dense square matrices

- Serial complexity,  $O(N^3)$

- block operations are useful

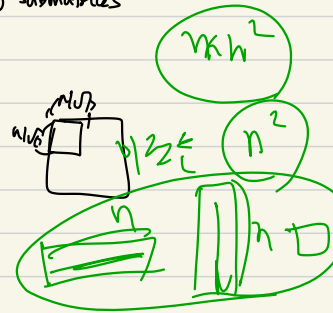


- $3^3$  matrix multiplication
- each using  $(N/4) \times (N/4)$  submatrices

## Block distribution

- decompose  $A$  and  $B$  matrices into  $P$  blocks

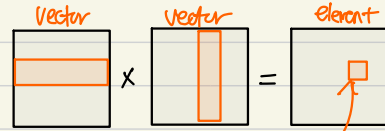
- each block size =  $(N/P) \times (N/P)$



- Process  $P_{ij}$

- stores  $A_{ij}, B_{ij}$

- computes  $C_{ij}$  of result matrix



## Computation

- all-to-all broadcast for row data

- Cost =  $t_s \cdot \log \sqrt{P} + t_w (N^2/P) (\sqrt{P}-1)$  (hypercube)

- all-to-all broadcast for column data

- Cost =  $t_s \cdot \log \sqrt{P} + t_w (N^2/P) (\sqrt{P}-1)$  (hypercube)

- compute local

- Cost =  $\sqrt{P} * (N/\sqrt{P})^3 = N^3/P$

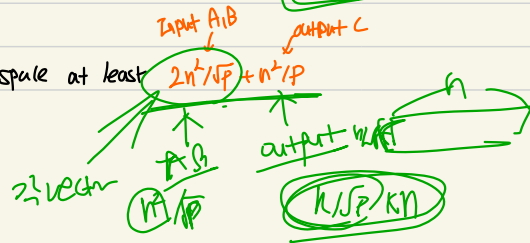
$$\sqrt{P} * (N/\sqrt{P})^3 = N^3/P$$

## Total parallel run time

- $T_p = \frac{N^3}{P} + 2t_s \log \sqrt{P} + 2t_w (N^2/P) (\sqrt{P}-1)$

## Memory overhead

- each node needs memory space at least



# Matrix-Matrix Multiplication (Cont.)

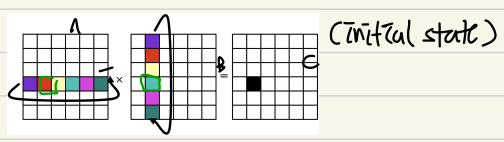
Memory overhead  $\approx 2n^2/p$

## Cannon's Algorithm

: memory-efficient version of MM algorithm

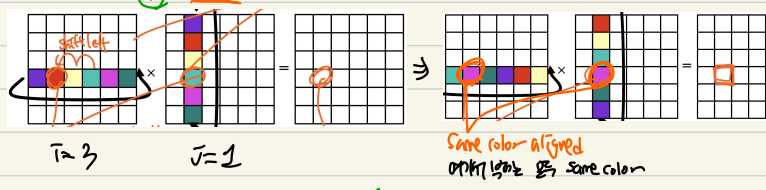
- Circulate the blocks of A and B with each process multiplying local submatrices at each step

Steps



① perform alignment until same color on target element

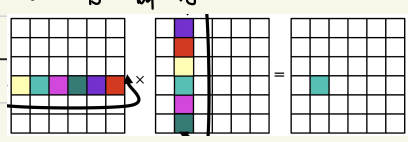
- Shift  $A_{ij}$  left by  $i$
- Shift  $B_{ij}$  up by  $j$



- Cost : Maximum Shift amount  $\sqrt{p}-1$  blocks  
each shift communication  $t_s + t_w n^2/p$  (block)

② perform local block multiplication and add to partial result

- Shift  $A_{ij}$  left by 1
  - Shift  $B_{ij}$  up by 1
- same shift amount ... total 6 shift needed



→ Multiplication step =  $\sqrt{p}$

- Cost :

- Communication :  $\sqrt{p}$  shifts and each shift takes  $t_s + t_w n^2/p$
- Computation :  $n^3/p$

alignment  $n * (\sqrt{p} * n^2/p)$   $\approx 3n^3/p$

$2(\sqrt{p}-1)(t_s + t_w n^2/p) + 2\sqrt{p}(t_s + t_w n^2/p) + n^3/p \approx 3n^3/p$

works more than serial but memory efficient

\* Cannon's algo  
local computation  $\approx$  partial reduction  $\approx$  discarding  $\approx$  local  
 $\approx (n^2/p)$  memory overhead  
 $\approx 2n^2/p$   
1-for input (temp)  
1-for output

Parallel Time

✓ 100%

# Common

10/11/2023

# 3D Matrix Multiplication

Algorithm	Comm	Compute	Memory per processor
Original MM	$O(n^2/\sqrt{p})$	$\theta(n^3/p)$	$\cong 2n^2/p^{1/2} + n^2/p$
Cannon's MM	$O(n^2/\sqrt{p})$	$\theta(n^3/p)$	$\cong 3n^2/p$
3D MM	$O(n^2/p^{2/3} \log p)$	$\theta(n^3/p)$	$\cong 3n^2/p^{2/3}$

[Input 2M]  
[Output m]

$M/\sqrt{p} \times M/\sqrt{p} \neq n^2/p$