

Lesson 2-4 Distributed Matrix Multiply

Matrix Multiply: Basic Definitions

$C \leftarrow C + A \cdot B \rightarrow$ This is the dot product of row A and column B and accumulating the sum into the output.

The Matrix Multiply as pseudo code:

```
for i ← 1 to m do
  for j ← 1 to n do
    for l ← i to k do
      C[i,j] ← C[i,l] + A[i,l] . B[l,j]
```

The Time to complete the algorithm is:

$T.(m,n,k) = O(mnk) \rightarrow T.(m,n,k) = O(n^3)$ when $m=n=k$

The Matrix Multiply as Parallel pseudo code:

```
parfor i ← 1 to m do
  parfor j ← 1 to n do
    for l ← i to k do
      C[i,j] ← C[i,l] + A[i,l] . B[l,j]
```

This means each 'row' and 'column' could actually represent a sub matrix.

The third loop is a reduction.

The Matrix Multiply as Parallel pseudo code:

```
parfor i ← 1 to m do
  parfor j ← 1 to n do
    let T[1:k] = temp array
    parfor l ← 1 to k do
      T[l] ← A[i,l] . B[l,k]
    C[i,j] ← C[i,j] + reduce(T[:])
```

$W(n) = O(n^3)$

$D(n) = O(\log n)$

A Geometric View

Using a cube - the rows and columns can be projected onto the cube. The three matrices are areas on the x, y, z planes. If the the three projections intersect - the matrices can be multiplied.

The resulting volume is the set of multiplications that need to be done.

According to Loomis and Witney: The volume of I is $||I|| \leq \sqrt{|sA| |sB| |sC|}$

If there is no intersection between the surface projections - there will be no multiplications.

1D Algorithms

How should the Matrix multiply be accomplished on a distributed machine?

Using Block Row Distribution : this means each node gets n/P rows. (assume the matrices are square and that n is divisible by P).

What is the most work that can be done by a node given the above data distribution?

Since a, b, c have to be multiplied - one of them has to be shifted, convention says 'B' is shifted.

1-D Block Row Distribution Pseudo Code

let $A'[1:n/P][1:n]$ = local part of A
 B', C' = same for B, C

let $B''[1:n/P][1:n]$ = temp storage

let $r_{next} \leftarrow (RANK + 1) \bmod P$

$r_{prev} \leftarrow (RANK + P - 1) \bmod P$

for $L \leftarrow 0$ to $P-1$ do

$C'[:,L] += A'[:,L] \cdot B'[:,L]$ (...L... is a placeholder for the local indices)

 sendAsync ($B' \rightarrow r_{next}$) (send the local buffer to the next processor)

 recvAsync ($B'' \leftarrow r_{prev}$) (receive from the previous processor)

 wait(*) (wait for send and recv to complete)

 swap(B', B'') (swap the receive buffer with the compute buffer)

1D Algorithm Cost

The cost of the algorithm is ...

$\tau = \text{time per "flop"}$ (flop means ...1 floating point multiply or add)

Total time is $T_{comp}(n,P) = 2\tau n^3 / P$

How much time is spent on communication?

B' is the only data communicated. It's size is n/P words by n columns - so n^2/P words.

There are P sends that have to be paid for.

So the total cost of communication is : $\alpha P + \beta n^2$

1D Algorithm Cost Part 2

The running time of the algorithm is flops + communication time.

$$T_{1D}(n;P) = 2\tau n^3 + \alpha P + \beta n^2$$

Getting more speed from the algorithm:

Here is the pseudo code, rearranging the code can improve the speed.

```
let A'[1: n/P] [1:n] = local part of A
    B', C' = same for B, C
```

```
let B''[1:n/P][1:n] = temp storage
let rnext ← (RANK + 1) mod P
    rprev ← (RANK + P -1) mod P
for L ← 0 to P-1 do
    C'[:,:] += A'[:,...L...] . B'[:,...L...][:]
    sendAsync (B' → rnext)
    recvAsync (B'' ← rprev)
    wait(*)
    swap(B', B'')
```

Rearranged pseudo code:

```
let A'[1: n/P] [1:n] = local part of A
    B', C' = same for B, C
```

```
let B''[1:n/P][1:n] = temp storage
let rnext ← (RANK + 1) mod P
    rprev ← (RANK + P -1) mod P
for L ← 0 to P-1 do
    sendAsync (B' → rnext)
    recvAsync (B'' ← rprev)
    C'[:,:] += A'[:,...L...] . B'[:,...L...][:]
    wait(*)
    swap(B', B'')
```

The communication is now done BEFORE the computations - gaining a little speed in the process.

This could improve running time by a factor of 2. We can see this as:

$$T_{1D, overlap}(n; P) = \max(2\tau n^3/P, \alpha P + \beta n^2)$$

This uses the fact that : $a + b \leq 2 * \max(a, b)$

Efficiency and the 1D Algorithm

Recall the running time: $T_{1D, overlap}(n; P) = \max(2\tau n^3/P, \alpha P + \beta n^2)$

Speedup: $S_{1D}(n; P) \equiv T_*(n) / T_{1D}(n; P) = P / \max(1, 1/2 * \alpha/\tau * P^2/n^3 + 1/2 * \beta/\tau * P/n)$

= $\theta(P)$

Parallel Efficiency = Speedup / P = E(n; P)

A parallel system is efficient if its parallel efficiency is constant. This occurs when:

$$n = \Omega(P)$$

According to the equation ... if you double the number of nodes - you have to double the dimension of the problem. But doubling the size of the dimension (since matrices are involved) will lead to quadruple the size of the matrices. The number of flops will increase by a factor of 8.

If you can't (or don't) double the dimensions you will see diminishing returns in increasing the parallelism.

Isoefficiency Function is: $n = \Omega(P)$ → the value of P that n must satisfy to have constant parallel efficiency.

Temporary Storage: the B'' needs temporary storage. The 3 matrices and 1 temp. matrix.

$$\text{Temporary Storage} \rightarrow M(n; P) = (3 + 1) n/P * n = 4 n^2/P$$

IsoEfficiency

$$E(n;P) = S(n;P)/P = T_*(n) / P * T(n;P)$$

$$\text{Parallel Cost} = P * T(n;P) = 1/((1 + P/\tau) \log P + (\alpha/\tau) (\log P)/n)$$

$\tau = \text{time per scalar add}$

$$T_{tree}(n;P) = \tau n \log P + \alpha \log P + \beta n \log P$$

Notice: there is no setting of n that will make E(n;P) a constant.

$1/((1 + P/\tau) \log P) \rightarrow$ goes to infinity as n increases.

A 2D Algorithm SUMMA

Begins with a 2D distribution of the matrix operands. Each node is responsible for updating the part of the 'C' matrix that it owns.

The SUMMA algorithm begins by integrating strips of width and height l.

So:

```
for l ← 1 to n/S do
    broadcast(horizontal strip, (owner))
    broadcast(vertical strip, (owner))
```

The run time:

Assume: nxn matrices, $\sqrt{P} \times \sqrt{P}$ mesh, \sqrt{P} and s both divisible by n.

$$\begin{aligned} T_{SUMMA}(n;P,s) &= n/s * (2\tau * (n^2/s/P) + T_{net}(n;P,s)) \\ &= 2\tau n^3/P + T_{net}(n;P,s) \end{aligned}$$

SUMMA Communication Time

$$T_{net} = O(\alpha * n/s * \log P + \beta * n^2/\sqrt{P} * \log P) \rightarrow \text{Tree}$$

$$T_{net} = O(\alpha * n/s * P + \beta * n^2/\sqrt{P}) \rightarrow \text{Bucket}$$

Efficiency of 2D SUMMA

The 2D SUMMA is more scalable than the 1D Block.

Note: 's' is the width of the strip, it is also the tuning parameter of the algorithm.

$$n_{tree}(P) = \Omega(\sqrt{P} \log P)$$

$$n_{1D}(P) = \Omega(P)$$

$$n_{bucket}(P) = \Omega(P^{5/6})$$

The bucket is slightly worse than the tree, it trades a higher latency cost for a lower communication cost.

SUMMA Memory

The amount of memory needed for SUMMA is:

$$M_{SUMMA} = 3 * n^2/P + 2 * s * n/\sqrt{P}$$

$$> 4 * n^2/P \text{ when } s > 1/2 * n/\sqrt{P}$$

A smaller 's' increases latency time.

If s is at it's maximum value, the SUMMA algorithm might need 5 times \sqrt{P} amount of storage.

A Lower Bound on Communication

lower bound - the number of words a node MUST communicate.

each phase sends and receives exactly 'm' words.

$S_A, S_B, S_C \rightarrow$ the set of unique elements of each matrix seen in this phase.

$$\text{Max \# multiplies per phase} \leq \sqrt{|S_A| * |S_B| * |S_C|} \leq 2 * \sqrt{2} * M^{3/2}$$

$$L \geq \# \text{ full Phases} \geq \lfloor W / \text{max \# multiplies per phase} \rfloor_{\text{FLOOR}}$$

$$L \geq W / (2\sqrt{2} * M^{3/2})$$

$$\# \text{ words communicated by 1 node} \geq (\# \text{ full phases}) * M$$

$$\# \text{ words communicated by 1 node} = \Omega(n^2 / \sqrt{P})$$

A Lower Bound on Communication

$$T_{net}(n;P) = \Omega(\alpha * \sqrt{P} + \beta * n^2/\sqrt{P})$$

$$T_{SUMMA, net}(n; P, s) = \left\{ \begin{array}{l} \alpha * n/s * \log P + \beta * n^2/\sqrt{P} * \log P \quad (tree) \\ \alpha * n/s * \sqrt{P} + \beta * n^2/\sqrt{P} \quad (bucket) \end{array} \right\}$$

(assume: $1 \leq s \leq n/\sqrt{P}$)

$$T_{Lower}(n;P) = \Omega(\alpha\sqrt{P} + \beta * n/\sqrt{P}) \quad \text{assume: } M=\theta(n^2/P)$$

Cannon's