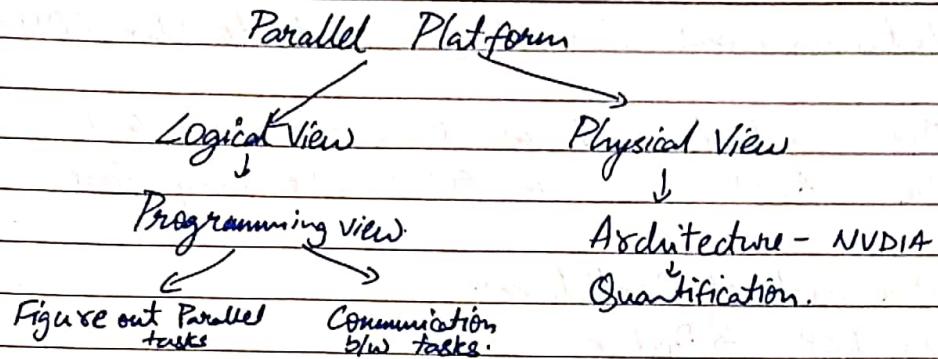


GPU



- Parallel Decomposition of Work: Division of tasks among the cores.
 - Resource Requirement: Cores required.
 - Mapping between resources & tasks.
(cores)
 - Communication b/w the cores for Synchronisation
due to dependency of a core's task on another core's task.
 - Latency - Time required for computation

Measured in seconds.
Bytes/sec.

Flops (Floating-point Operations) → Data Access (Majority of Latency)
(Latency)
 - Bandwidth - Bytes/s
(Memory throughput)
 - Compute Throughput - FLOPs/second
Focus on latency
 - Latency devices - Good in doing a single operation quickly
(CPU)
Focus on throughput 1 job - 1 min
10 jobs - 10 min
 - Throughput devices - Good in doing multiple jobs at a time
(GPU)
but not good for single op.
1 job - 1 min
10 jobs - 5 min.

→ Heterogeneous Computing: Used to Use CPU + GPU

→ The control unit in CPU's are very fast
The ALU in GPU's are very fast.

→ Three levels of Parallelism:

(i) ILP → Pipeline.

(ii) Thread → GPU → CUDA

(iii) Node → ~~CPU + GPU~~ Multiple GPUs

→ Multi-Core → CPU - 6 cores.

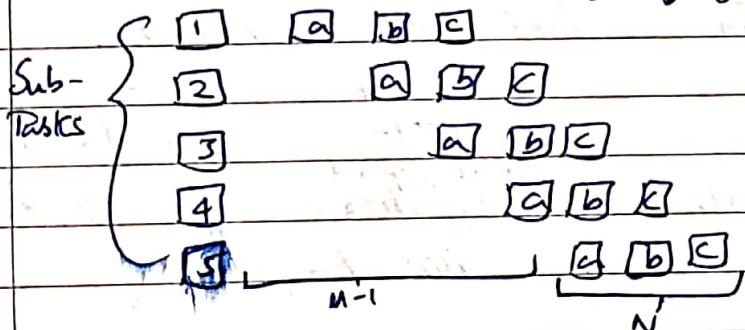
Many-Core → GPU - 100s of cores

→ ILP:

↳ Instruction Level Parallelism:

The pipeline divides a task into sub-tasks.

↓ ↓ ↓ ↓ ↓ ↓ → + (cycles).



Sub-tasks are distributed across functional units specialising in the sub-tasks & carry it in a single cycle.

If ~~the~~ task was carried out in serialised manner:

$$T_{\text{serial}} = m \times N \text{ cycles}$$

No. of subtasks

$$\text{In ILP, } T_{\text{ILP}} = (m-1) + N \text{ cycles.}$$

Since $(m-1) \ll N$, $T_{ILP} \approx N$

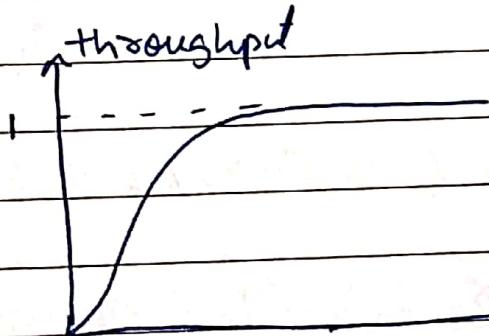
So Speed up = $\frac{mN}{N} \approx m$.

more

The higher the depth the better the speed up.

- Each sub-task is a thread.
- Wind up time = $m-1 \rightarrow$ No output for that time.
- All threads can also have the same task & there is no need for unique sub-tasks.
- ILP is not helpful when m is small comparable with N . Ex: tight loops \rightarrow 5 iterations.

→ Instruction throughput = $\frac{N}{N + m-1} = \frac{1}{1 + \frac{m-1}{N}}$



The ideal situation is one instruction per cycle.

So max. throughput is 1/st/cycle
and thus is the ILP-wall.

- Multiplicity is used to improve computing i.e. compute time ↓
Data access

↓
Multiple Memories & Multiple data paths

Multiple cores
CPU - 2, 4, 8...64 GIPUs - 1000

- In CPU memory allocation, location cannot be specified
GPU can be specified.

Hierarchy of Searching for data

→ Processing Unit

Registers

L1 Cache

L2 Cache

L3 Cache

RAM

Present on each cores

Shared by the cores

- If a cache is k times faster than the RAM
Re-use ratio is γ → temporal locality of data.

T_m - time reqd. to get data from RAM

$$T_c = \frac{T_m}{k}$$

Cache Latency

$$\text{Avg time to get data} = \gamma T_c + (1-\gamma) T_m$$

$$\text{Speed-up} = \frac{(\gamma T_c + (1-\gamma) T_m)}{T_m} = \frac{T_c (\gamma + (1-\gamma) k)}{k T_c}$$

$$= \frac{k T_c}{T_c (\gamma + (1-\gamma) k)} = \frac{k}{\gamma + (1-\gamma) k}$$

If $\gamma = 1$, Speed up = k .

If $\gamma = 0$, Speed up = 1 → no speedup.

% operator is much slower than + or - since the pipeline is much longer

11

→ Spatial Locality:-

Cache Lines: If $a[i]$ is requested, $a[0], a[1], \dots, a[n-1]$ are fetched

$$\rightarrow \text{FLOPS} = \frac{\text{No. of Cycles per Second} \times \text{No. of FLOPs per cycle} \times \text{No. of cores per processing unit}}{\text{Clock rate}}$$

$$\text{For a CPU, FLOPS} = \frac{2.5 \times 10^9}{2.5 \text{ GHz}} \times 4 \times 1 \\ = 10 \text{ Giga Flops.}$$

$$\text{For a GPU, FLOPS} = \underline{10^4} \text{ Giga Flops} = 10 \text{ TFlops.}$$

But for a GPU → Clock rate is $\approx 1 \text{ GHz}$

$$= \frac{10}{2.5} = \underline{4 \text{ TFlops.}}$$

Data rate is in GigaBytes per second & Flops is in Teraflops
So parallelism is required to match them.

→ Missed Clash - 1/8/19

→ Vector Triad → Low Level Benchmarking → Small standard codes to test ordered total is to ensure that no. of iterations are constant

Properties of Vector triad:

Load → 3

Store → 1

for $i=0 : R$

for $j=0 : N$

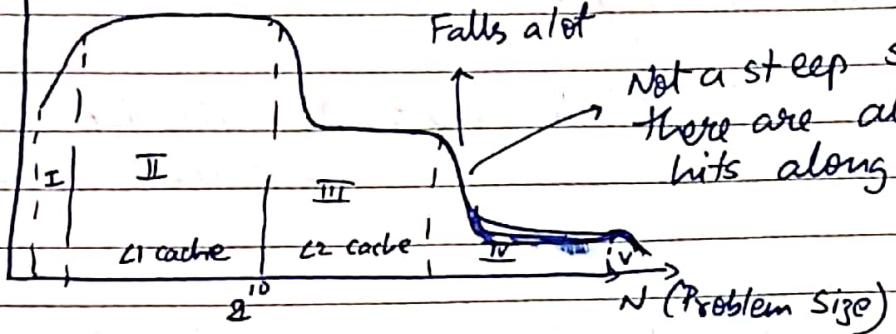
$$a[j] = b[j] + c[j] \times d[j]$$

Operations → 2

Spatial Locality & Temporal Locality

Throughput
(GFlops)

$$\text{Double Matrix} \\ A[i][j] = B[i][j] + C[i][j] * D[i][j]$$



- Initially pipeline does not offer any benefit because no. of operations are small.
- Bash Command: lscpu → Memory size
A, B, C, D values
- II ends at 2^{10} operations - $2^{10} \times 8 \times 4 = 2^{15}$ bytes → L1 cache size.
double
- L1 cache can be used by just one processor because it is present on each core.
- Bash Command: top : Processes Running
ps -a :
- III to IV drops a lot → maybe because L3 is shared.
Not entirely
But the drop would remain the same even if there is a single process & it is allowed to access unlimited memory

Memory Throughput

- L1 & L2 → Different latencies but same bandwidth
- L3 → Different latencies & bandwidth w.r.t L1 & L2

→ Describes an algorithm.

Computational Intensity

→ Code Balance \hookrightarrow Compute to Memory Access Ratio
(CMA Ratio). (Flops/Bytes)

For Vector triad:

i) If Store is considered $\Rightarrow \frac{2}{4} = \frac{1}{2}$

ii) " " " ignored $\Rightarrow \frac{2}{3}$.

Can be ignored if RUS has a lot more data access

GPU is suitable for higher CMA.

Computation Intensity = 1

Code Balance \hookrightarrow Bytes/Flops

→ Machine Balance = $\frac{\text{Theoretical Max. Memory Bandwidth}}{\text{Theoretical Max. Compute Performance}}$
(Bytes/Flops)

↳ Describes an

Architecture \rightarrow So fixed for a computer.

~~Lower Code Right now Code Balance is increasing~~
~~but Machine Balance is decreasing.~~
(current trend)

→ CPUs can be 10x faster for sequential code
GPUs can be 10x faster for parallel code.

→ Constraints on parallel programming \rightarrow Time
↳ Resources

↳ Accuracy \propto No. of Computations.

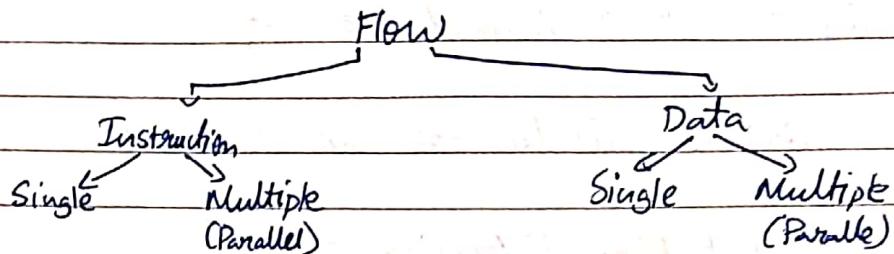
→ Memory system's consume a lot of power.

→ Vector processes $\rightarrow A[i] = B[i] + C[i]$

Instead of doing it index by index.
entire chunks are computed at a time

→ Von Neumann bottleneck → The increasing gap b/w the data access & processing architectures (Performance-wise)

→

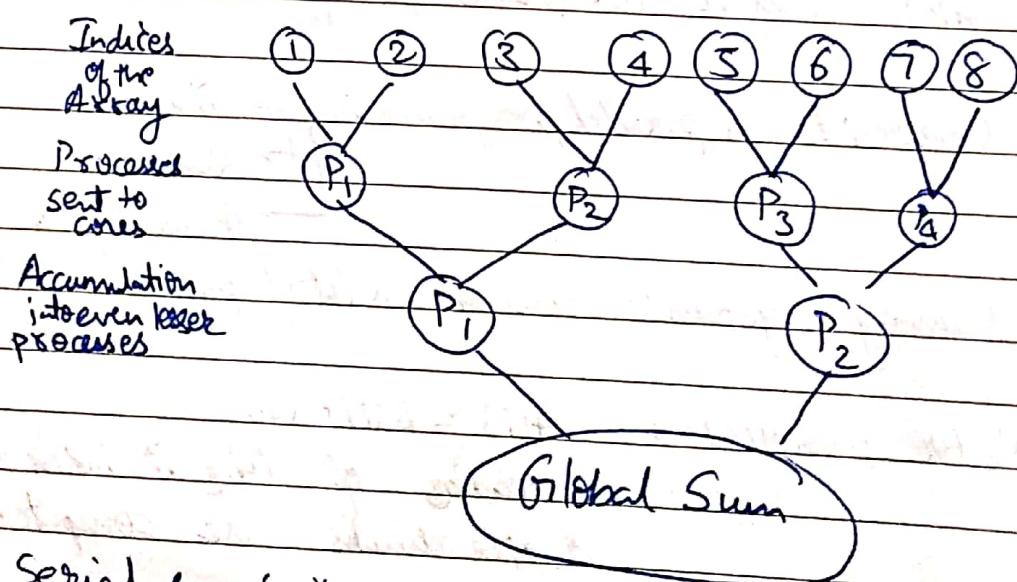


So Flynn's Taxonomy :	Single Instruction	Single Data
	Multiple	" "
	Single	Multi "
	Multi ..	" "

GPU programming is usually SI Multi-thread(Data) (SIMT) / (SIMD)

→ Dependency Analysis:-

Sum of nos. in an array :



Serial complexity $\rightarrow O(n)$
Parallel complexity $\rightarrow O(\log n)$

$$\text{Speedup} = \frac{N}{\log_2 N}$$

Such a speedup is possible for any associative operation.

For such a process distribution analysis we assume:

- * Memory Access time is neglected
- * No. of available cores is comparable to size of the problem
- * Each operation takes same time
- * Efficiency of each core is the same.

$$\rightarrow \text{Theoretical Maximum Speedup} = \frac{N}{N/P} = P.$$

But memory cache can be used to get a practical speedup $> P$. \hookrightarrow Since combined core cache \gg Single or serial cache

$$\rightarrow \text{Efficiency} = \frac{\text{Speedup}}{\text{Processes}}$$

Sum of nos. in Array (Contd.):

Serial & Parallel Work Complexity = $O(n)$.

Sometimes parallel work complexity $>$ serial. \rightarrow due to dependency
But time complexity may be lower.

\rightarrow Amdahl's Law \rightarrow Theoretical speed up given enough processor

$$\text{If } p=1, \text{ Speed up} = \frac{1}{1-1} = \frac{1}{0} \rightarrow \infty.$$

\hookrightarrow Only when ∞ processors are given

gprof --line a.out gmon.out
gcc -pg file

Speedup = $\frac{1}{\frac{P}{N} + S}$ → fraction of code that can be parallelized at that time \rightarrow No. of processes.
can be parallelised. So here S determines the scope of parallelization.

→ Dependency:

I $\text{for } i = 1 : N$

II $\text{for } j = 1 : M$

$$b[i,j] = b[i,0] + \dots + b[i,j-1].$$

Also
for d

Since there is a dependency on $j-1$ for j only loop I can be parallelised.

→ Flow dependence : Statement i precedes j & (True dependence) i computes a value j uses.

(i) $x = 1$

(ii) $y = x + 2$

→ Anti dependence : Statement i precedes j & i uses a value j computes.

(i) $y = x + 2$

(ii) $x = z - w$

→ Output dependence : .. " " " i computes a value j computes

(i) $z = 3 - w$

(ii) $y = z / 3$

→ Input dependence : .. " " " uses a value j uses

(i) $y = x + 2$

... - - -

→ Loop Independent Dependence:

Loop {

$$a[i] = b[i] + c[i]$$

$$d[i] = a[i]$$

Dependence distance is 0

No. of iterations over
which dependence exists

→ Loop Carried Dependence:

Loop {

$$a[i] = b[i] + c[i]$$

$$d[i] = a[i-1]$$

Dependence distance is 1.

Also have
flow
dependence

→ $\text{for } (i=0; i < N; i++)$ → Can be parallelised
 $\text{for } (j=0; j < N; j++) \{$ → Cannot be parallelised.
 $a[i][j+1] = a[i][j] + 1.$

→

$\text{for}(i=0 : N)$

$$a[i-1] = b[i]$$

$\text{for}(j=0 : n)$

$$a[i] = a[j]$$

can't be parallelized

Loop Fusion

$\text{for}(i=0 : n)$

$$a[i-1] = b[i]$$

$$c[j] = a[i]$$

Cannot be parallelized.

Better for ~~steps~~
Serial running due
to locality of reference

Loop Distribution

→ CUDA - Compute Unified Device Arch.

→ GPU Kernels - Functions · Usually SIMD

In a program some part is executed on CPU & GPU
The parts going to GPU - Kernel.

→ GPU & CPU have their own RAM & memory.

→ Stack vs. Heap

* Scopes are different. Stack - once out scope, memory is

* Stacksize is static. Heap is dynamically allocated cleared.

* Size of Heap is much larger than the size of stacks.

* Bookkeeping mechanisms are different.

* Programmer manages heap but OS manages OS

→ CPU manages the ~~normal~~ CUDA world. P ~~operations~~
Since CPU's archi. supports control.

Taken

→ Threads → Cores

↓
Cores

Blocks → Streaming Symmetric Multiprocessor (SM)

↓

Grids → Device (GPU)

- A thread is a lightweight process
- A process has * memory - Stack, Heap, Code, data
 - * Process Control Block - Stores its
 - ↳ Stores PID
 - ↳ Stores Open files - to allow context switch
 - ↳ Stores Register values
- A warp is 32 threads currently - can change depending
 - * 1 SM runs a warp
 - on the capabilities of the GPUs in the network
- cudaMemcpy → a form of direction is required, because
 - ~~and~~ the same address can exist in both CPU & GPU → so details regarding direction is required.
 - and blockid, blockDim
- threadIdx is a struct generated
 - can't be accessed only inside the kernel.
 - The x, y, z elements show the id in the axis direction.
- Registers are very valuable in GPU & therefore extra variables must be minimized.
- The code might look like:
 - kernel → Run by GPU waiting for Making the CPU wait while the kernel runs wastes time
 - SerialCode → Run by CPU
 - Kernel
 - SerialCode

When operations are asynchronous, other ~~other~~ operations can be run concurrently.

→ Scheduling threads takes 100s - 1000s of clock cycles
in CPUs

takes 2-3 clock cycles in GPUs.

because there is no-preemption in GPU scheduling

→ --device-- → executed on device, callable from device

--global-- → " " " " , " " host

--host-- → " " host , " " host

(can be used together i.e a function can be both)

Ques: 2x

→ 0 1 2 3 4 5

→ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11

→ Constraints:

* No. of threads per block

* No. of blocks per SM

Such constraints can be known with `cudaGetDevice`
function.

Architectures { Fermi, Tesla, Kepler Major Minor Improvements
 ↑ to arch. }

→ Compute Capability → Hardware Ex: 3.1, 1.6

CUDA Toolkit Version → Software Ex: 10.1

Ques:

$$\rightarrow id = 2 \times (\text{Blockid} \cdot x + \text{Blockdim} \cdot x + \text{threadid} \cdot x)$$

ie

$$\rightarrow \text{Ques 2: } id = 2(\text{Blockdim} \cdot x \times \text{BlockID}) + \text{threadid} \cdot x$$

\rightarrow Image Processing.

↳ Memory is flat & not in 2-D.

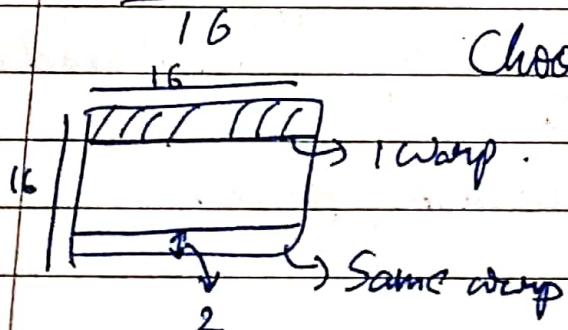
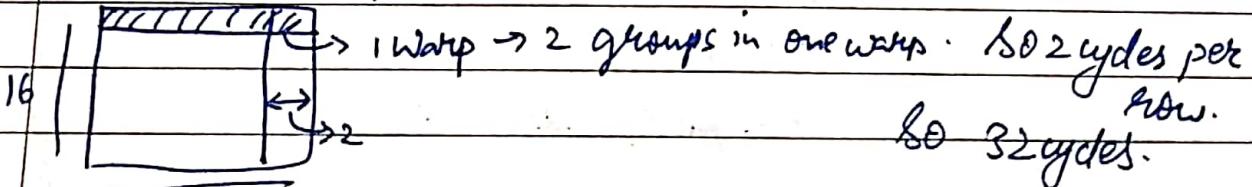
\rightarrow Latency hiding scheduling - Allocate any block which whose requested resources are available.

\rightarrow Since the Control Unit (CU) is small in the GPU

When all the threads have the same instructions then all of them ^(in a block) are executed concurrently - SIMD

Even if If there are n groups of threads in a block which have n different instructions then the n groups will run in different phases serially. Possible
 ↳ n Control Divergence.
 (can be less than n)

\rightarrow A warp is executed at the same time (32 threads).



Choosing different block size can help too.

\rightarrow 16 cycles.

- So the scheduling happens for warps & not threads unlike CPU.
- Stencils - The ~~is~~ data access organisation of a thread.
- Matrix multiplication of matrices M & N → P

$$P_{i,j} = M[i,:] \cdot N[:,j]$$

↳ dot product

Each thread is mapped to an i, j in P

- -- syncthreads() will ensure that ~~all~~^{no} threads will ~~not~~ cross the function until all threads have reached the function.
- Must be executed by all threads of the block and it must be the same function call.
i.e

if

x x x

-- syncthreads() - (1)

else

x x x

-- syncthreads() - (2)

- (1) A thread cannot execute (1) & another thread in the same block executes - (2)
- If this happens waiting forever.

2 times the bandwidth of global memory & also consumes an order of energy lesser

Very fast

Local Memory - R/W per thread

Registers - .. " " → A variable stored here will not have a copy on global

Shared Memory - R/W per block.

Takes time Global Memory - R/W per grid & R/W from host

Constant Memory - R per grid & R/W from host

Needs a load instruction

→ Every variable has a scope & lifetime

↓
if a thread or

If it is erased after

a block of threads
can access it

↓
the kernel or after entire process

declare in kernel

outside the func.

Lifetime Kernel

but Shared is faster

→ Automatic Variables (not Arrays)

Memory

Register

Thread

Scope

Thread

Lifetime Kernel

→ Automatic Arrays Local Thread Kernel

Declared in kernel or devicefunc.

→ --device-- --shared-- int x Shared Block Kernel
Optional

→ --device-- int x Global Grid Application

Declared outside function

→ --device-- --constant-- int x Constant Grid Application
Optional

Tiling:-

→ In matrix multiplication, $t_{0,0}$ and $t_{0,1}$ both access the 0th row of M. So if $t_{0,0}, t_{0,1}, t_{1,0}, t_{1,1}$ are in a single block & all corresponding rows are brought to shared memory - no. of accesses to global halves. If size of tile is $N \times N$ - " " " " reduces time.

One file

→ It is better to have 2 threads that access the same memory at a time to execute at the simultaneously

global to Shared is one operation

— / —

→ Reduction Operator

Eg: Addition, Max, Min of an array

- Associative
 - Commutative
- To use this pattern, the algo must be both.

We have a running value that is updated in each iter.

→ 

Each thread ~~can't~~ deals with 2 elements & so on
- $\log_2 n$ steps.

In global memory, data access is 2 & computations is 1
So CGMA = $\frac{1}{2}$. (very bad)

So first, transfer data to shared memory - using N threads.
- CGMA = 1

and then work by divide & conquer.

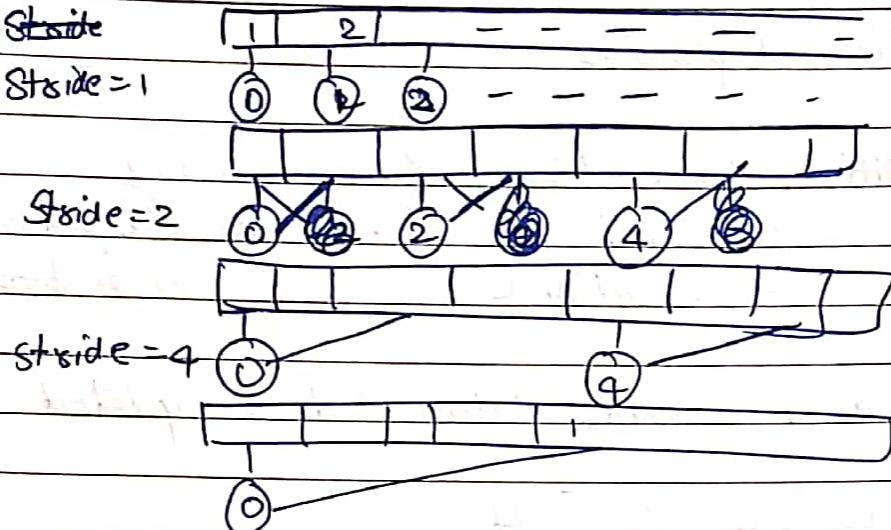
~~divide & conquer~~

Total No. of operations is $(N-1)$

So work complexity $O(N)$

Step complexity $O(\log N)$

↳ In all parallel algos benefit is in terms of Step complexity



→ `for (stride = 1; stride < blockdim.x; stride *= 2)`

- `if (tid % stride == 0)`

 Active

else

 deactive.

Synchronise.

→ Implementation - Stride Update.

- Activate/Deactivate.

- Synchronisation

→ Evaluation - Shared Memory

- CGM4

- Code Divergence.

) Code Divergence is horrible

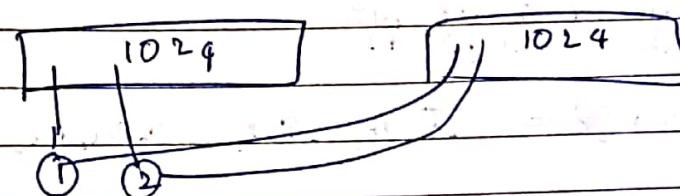
Replace with `if (tid < blockdim/stride)`

 Activate

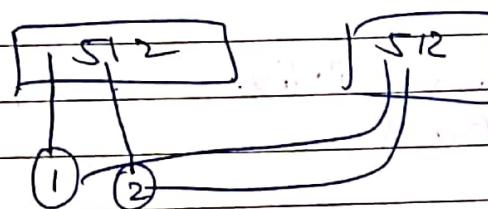
 else deactivate.

Then consecutive threads will either be activated or de-activated

If $N = 2048$ then divide into 2 sections.



Then $N = 1024$



→ Allocation - Done for a sm. by blocks.

Even if resources are not available for a single block - the block is not allocated to the sm.

Scheduling in terms of warps

) for (stride = blockdim, stride >= 1, $\frac{tr}{side} / = 2$)

if (hid C stride)
active.

→ ? Why load into shared - copy value in shared after first operation.

→ SCAN :-

I/P 1 2 3 4 5

O/P 1 1+2 3+2 5+3 6+4 (12+5)
(3) (15) (18) (12) 17

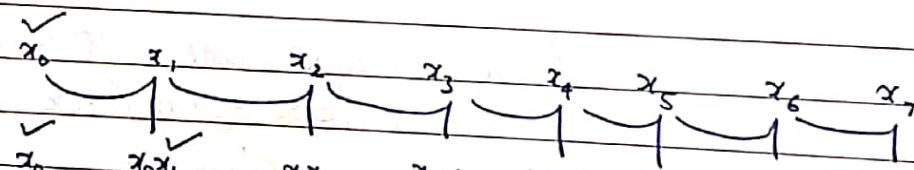
$$\text{scan}(A) = [I, a_0, a_0 \oplus a_1, a_0 \oplus a_1 \oplus a_2, \dots, a_0 \oplus a_1 \oplus a_2 \oplus \dots \oplus a_{n-1}]$$

Prefix scan if \oplus is addition
(Inclusive)

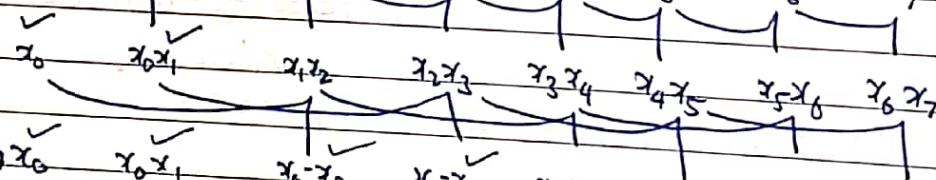
Inclusive begins with first element & last element is the sum
of all elements

Exclusive
Inclusive sum begins with 0 & last element is the sum
until the $(n-1)$ th element.

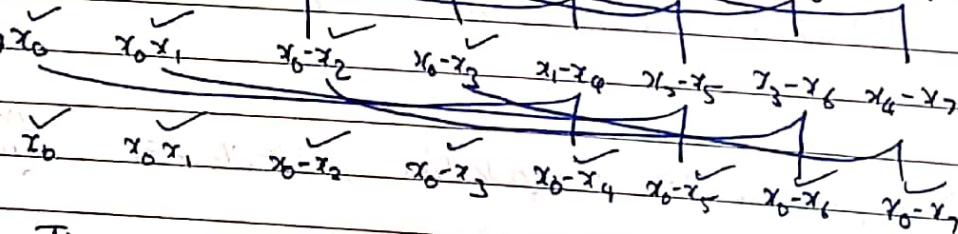
Stride=1



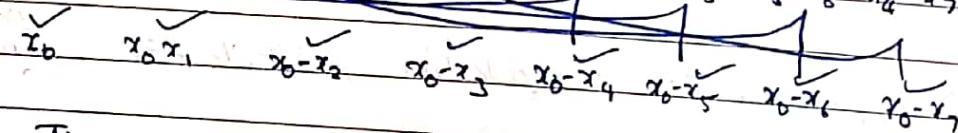
Stride=2



Stride=3



Stride=8



Time Complexity - $\log_2 N$

Work Complexity - $(N-1) + (N-2) + (N-4) + \dots - N \frac{\log N}{2} - (N-1)$ $\log_2 N$ time

$$- N \left(\log_2 \frac{N}{2} \right) - 1$$

So we have our work complexity being $\log_2 N$ times
that of the serial complexity

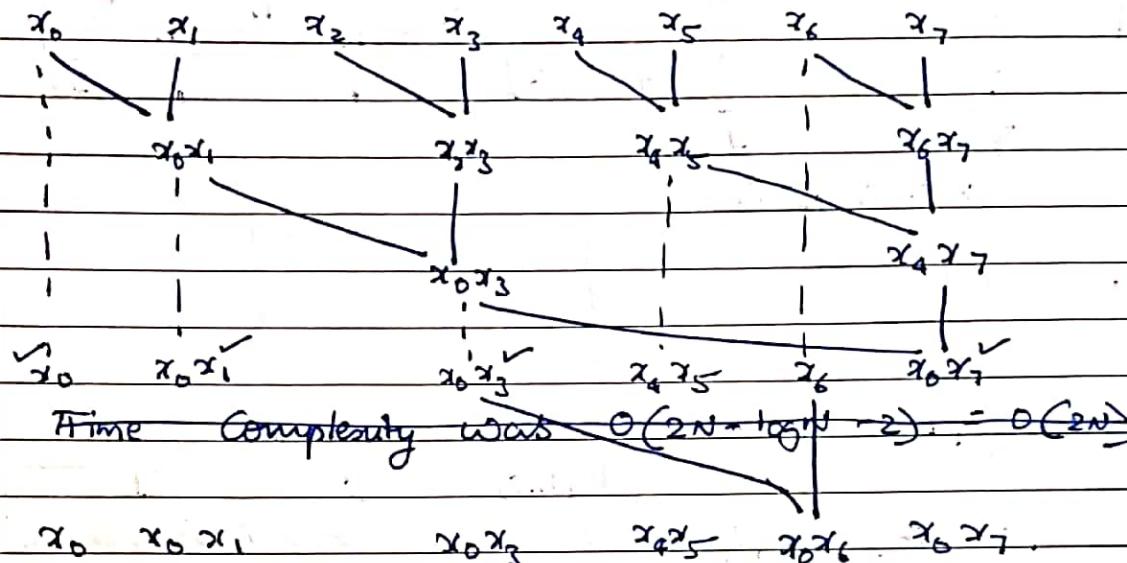
~~888888~~

$$N-1 + N-1 - \log N$$

$$2N - \log N - 2$$

— / —

So we need atleast $\log N$ threads to achieve the speed-up equal to serial code.



Time Complexity is $O(2N - \log N - 2) = O(2N)$

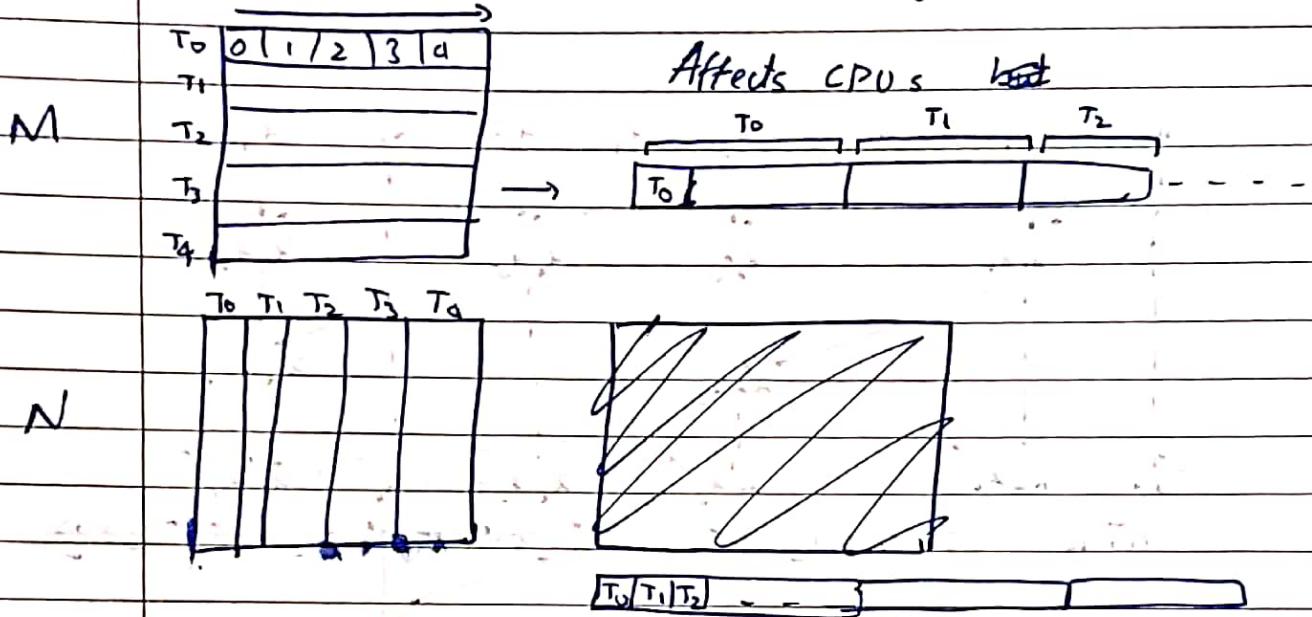
We divide them into blocks \Rightarrow block 1 depends on block 0.
block 2 " " " block 0 & 1

So we do a scan of the last values in the blocks.
then add these to all values

→ ↴

'False' → Coalesced

→ Coalescing - the direction we move in is the direction in which memory is arranged



At $t_0, T_0, T_1, \dots, T_n$ all access their memories
so it is easier for the warp to pack the memories in the N for t_0 than for M for t_0 .

Convolution:

○ Element Reduction Called embarrassingly parallel algo

→ Bigger the control divergence if mask size ↑

○ Control Reduction without shared memory issue

→ Constant Memory is the cache

Quiz:

0	1	2	0	4	5	0	0	8	9
(0)	1	1	0	1	1	0	0	1	1
	1	2	2	3	4	4	5	6	

Scan

- Prototype Mask/filter can be stored in Constant memory
- Each thread makes $(2n+1)$ global memory accesses
Filter size. & $(2n+1)$ multiplication
& $(2n+1)$ additions.
So CGMA = $\frac{4n+1}{2n+1} \approx \frac{4n}{2n} = 2$ (not good)

→ To increase CGMA, we use shared memory but in block if tid get tid index's element then we need n elements to the left & right of the block.



↳ Control divergence

$$CGMA = \frac{\text{block dim} \times 4n}{\text{block dim} + 2n} \approx \frac{\text{block dim} \times 4n}{\text{block dim}} = 4n$$

\propto Size of Mask

Matrix Compression:

→ Compressed Sparse Row:

1	0	0	2
0	0	0	0
2	3	0	0
0	0	1	0
0	0	0	2

Data = {1, 2, 2, 3, 1, 2}

Col = {(0, 3), (0, 1), (2, 1), (3)}

Row = {0, 2, 2, 4, 5, 6}, → No. of elements in the row.

→ No. of elements

Starting index of the row

~~a < c~~ ~~d < c~~

~~b < c~~ ~~a < c~~
~~b < c~~ ~~d & b is at~~
→ Histogram :-

— / —

~~Due~~ Due to data coalescence (I) is better than (II)

I) $[T_0 \cap T_1 \cap T_2] - - - [T_3 \cap T_4] - - - [T_5 \cap T_6] - - - [T_7 \cap T_8] - - -$



II) $[T_0] - - - [T_1] - - - [T_2] - - - [T_3] - - - [T_4] - - -$

In the case of matrix mult. $M * N \rightarrow M$ is coalesced ^{memory access}
 $& N$ is partitioned ^{memory access}

→ If all threads are doing the same thing - maybe with different memory locations then there is no control divergence

→ When we have $T_0 - T_4$, all having 'p' as their characters → all access same bit
"Race Condition"

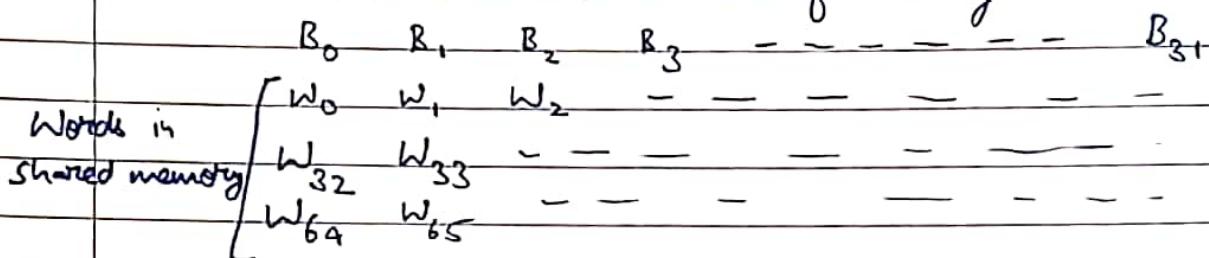
→ Bank Conflict :

Occurs when:

- * A warp is sharing/accessing data in a given pattern
- * Word = 4 bytes

Word - 4 bytes

Each threads has a bank of memory:



If shared memory is 8000 Bytes - 2000 words.

⇒ The memory is $w_0, w_1, w_2, \dots, w_{199}$ per clock cycle

Bank bandwidth - 32 bits/bank (fixed)

So if two threads access a bank at

Individual banks for each thread. the sometime, we will have latency.

So if $A[\text{threadIdx} \times 0]$ → Bank conflict at B_0

$\leftarrow A[\text{threadIdx} \times 1] \rightarrow$ No bank conflict

$A[\text{threadIdx} \times 2] \rightarrow$ Bank conflicts at every even

$\Rightarrow T_0 \& T_16$ access $w_0 \& w_{32} \rightarrow B_0$ by 2 threads bank

Broadcast

Multicast

(slower than Broadcast)

Two-way bank conflict

$A[\text{threadIdx} \times 3]$ since we only have 32 threads there is no bank conflict.

If multiple threads use the same memory (not bank)

Ex: $A[TID \times 0] \rightarrow$ Broadcast

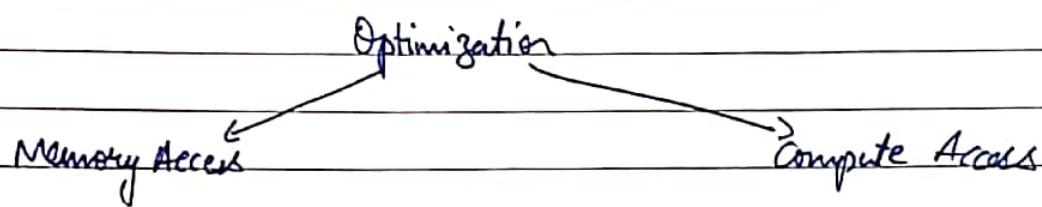
If divided into 2 or 3 memories

Ex: $A[\text{if } \dots] \rightarrow A[TID \cdot 1.2] \rightarrow$ Multicast

Broadcast is faster than Multicast.

Optimising Parallel Reduction in CUDA

- Even if only 1 byte is required, 1 full word is read.
→ If 8 bytes ~~are~~ array, then even $A[TID]$ has bank conflict.
- It is not necessary that all the words in 8 bytes or 12 bytes DS are being accessed at the same time.



Algorithm can be Memory Bound Compute Bound.

Reduction - CGMA = $\frac{1}{1}$ proportion Memory bound.

For memory optimization, parameters - memory bandwidth (GB/s)

compute optimization

- Compute throughput (GFlops/s)

- $\text{-- sync threads}()$ synchronises blocks in a block but not among blocks — kernel exit & launch must be used for sync b/w blocks.
- kernel launch has negligible overhead unlike thread launch in CPU.

#1

- Interleaved Addressing - ~~Bank conflict~~
- Too many deactivated threads
 - I. is a very slow operation
 - Code divergence

#2

- 'Solution' by multiplying tid with $2^5 \rightarrow$ code divergence
2. problem solved.
Bank conflict increased.