

13-01-2020

Sri Ramajayan

CS6023 : GPU PROGRAMMING

We think in a sequential manner, inherently. But we would like to learn how to write parallel programs, because sequential programs are slow in practice. ↳ to gain performance.

1976 → 1993 → 1997. Architectural improvements, clock Automatic improvement in performance.

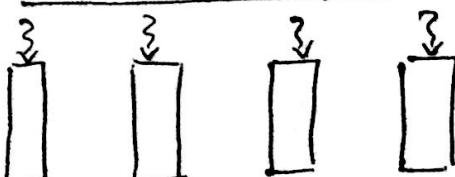
Trend could not be sustained. - cannot further ↑ clock frequency.
heat problems due to extreme power consumption of core.

what to do? - Multi cores.

No. of transistors on a chip ↑ clock speed, frequency, power ×.

Parallel Platforms. [cat /proc/cpuinfo → no. of cores].

1. Shared memory systems (multicore) $\sim 10^5$ threads - parallelly run on diff cores
RAM - shared memory among cores.



L₁, L₂ cache - local to individual cores, not shared

cores.



RAM/memory - shared among all cores.
(L₃) - shared across cores.

Why these many levels of cache?

→ Size v/s cost.

→ spatial, temporal locality of reference.

→ Speed matching. CPU is much much faster than RAM.

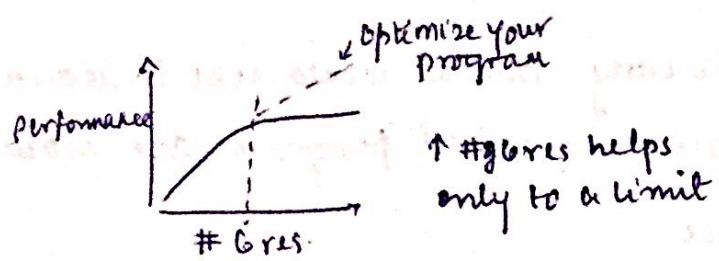
Bridging the speed gap b/w RAM, CPU.

2. Distributed systems (cluster). eg. cloud computing,

many diff. Computers connected by a network.
(nodes)

sequential ↳ parallel

3. Graphics Processing Unit (many cores): ~1000s.



4. Field Programmable Gate Arrays (FPGAs) - Configurable after manufacture
→ Specific to a particular application → create a chip

5. Application Specific Integrated Circuits (ASICs)

GPU performance - much better than CPU performance.

EVALUATION:

- 5 assignments ($7 + 7 + 13 + 13 + 20$)
 - Mid Sem (20)
 - End Sem (20)
- } Theoretical.

Amazon cloud.
Google colab. } GPU - get.

Cuda Compiler.

nvcc hello.cu

dkernel <<1,1>> () ; Kernel launch.

runs as single thread, on GPU

function kernel

CPU and GPU run asynchronously.

cudaThreadSynchronize();

HW

nvcc
cuda version
device query

Ref TB. Hwu, Kirk - book is slightly complicated initially.

Slides are available.

14-1-2020.

- whereis nvcc
- nvcc --version
- deviceQuery → CUDA SDK.
↓
properties of GPU device on our machine.

* Kernel launch in CUDA. asynchronous.

```
cudaDeviceSynchronize();  
  
dkernel <<< 1, 32 >>>();  
~~~~~  
product = 32 ⇒ prints 32 times.
```

Parallel Programming Concepts

1. Process: a.out, notepad, chrome - program in execution. ps, top.
2. Thread: lightweight process; a process gets divided into multiple threads or a part of it.
3. Operating System: windows, linux, Android.
4. Hardware: cache, memory, cores
5. Cores: threads run on cores; a thread may jump from one core to another.

Each thread has its own stack. (LIFO)

~~Global memory~~,
Heap ⇒ shared across threads.

1 core - 1 thread at a time
we can launch more threads than # of cores.

But it may not be a good idea to launch too many threads on CPU.

→ Context switch overhead.

On GPU, it's easy to launch large no. of threads.

Ideal to launch no. of threads = no. of cores to execute parallel programs.

- a) Why would a thread jump from one core to another? "sleep"
 $n = a[i]$ $\stackrel{\text{ld add of } a, R^1}{=}$ dependence $\stackrel{\text{load} \Rightarrow \text{wait}}{\text{load}}$ \Rightarrow wasting cycles.
 $=$ $\stackrel{\text{add } v \dots, R^2}{=}$ $\stackrel{\text{ld } [R^1], R^2}{=}$
- - $\stackrel{\text{ld } R_2, x}{=}$

Pick up some other core's thread which is "ready"

- Thread migration improves utilization of resources.
- Programmatically, we can restrict the threads from switching b/w cores. This is called "thread pinning."
- * Context of thousands of threads is stored in a dedicated set of registers? in GPUs?

32 *warp - warp based execution.

↓ warp scheduling - need not be executed sequentially.
set of threads. Need to write programs so that warp scheduling does not affect final result.

warp size - fixed for the time being (32) - NVIDIA GPUs

Q) convert sequential C code to a parallel CUDA Cde.

```
#include <stdio.h>
#define N 100
int main() {
    int i;
    for (i=0; i<N; ++i)
        printf("x.%d\n", i*i);
    return 0;
}
```

```
#include <stdio.h>
#include <cuda.h>
#define N 100
__global__ void dkernel() {
    printf("x.%d\n", threadIdx.x *
```

```
int main() {
    dkernel<<<1, N>>>();
    cudaDeviceSynchronize();
    return 0;
}
```

output
order would be
jumbled up

| Variables declared
| inside a kernel are
| local to that thread.

Thread ID \Rightarrow ??

threadIdx.x

dkernel<<<1, 10>>>(); thread Ids \Rightarrow 0...9
dkernel<<<1, 20>>>(); thread Ids \Rightarrow 0...19

a) Convert to CUDA

```
#include <stdio.h>
#define N 100
int main() {
    int a[N], i;
    for (i = 0; i < N; ++i)
        a[i] = i * i;
    return 0;
}
```

How to do it? What to do?

```
#include <stdio.h>
#include <cuda.h>
#define N 100
int a[N]; // CPU memory; cannot be accessed from GPU
X --global -- void func(int*) {
    a[threadIdx.x] = threadIdx.x * blockIdx.x;
}
int main() {
    func(&a, N);
    cudaDeviceSynchronize();
    return 0;
}
```

How to do it? What to do?

```
#include <stdio.h>
#include <cuda.h>
#define N 100
--global -- void func(int*a) {
    a[threadIdx.x] = threadIdx.x * blockIdx.x;
}
int main() {
    int a[N], *da;
    int i;
    cudaMalloc(&da, N * sizeof(int));
    cudaMemcpy(da, a, N * sizeof(int), cudaMemcpyHostToDevice);
    for (i = 0; i < N; i++) {
        printf("%d\n", a[i]);
    }
    return 0;
}
```

allocates mem on GPU
not essential, but good practice.
→ cudaMemcpyDeviceToHost);
similarly Host to device.

17-1-2020

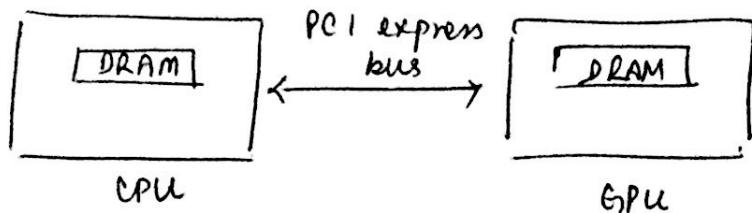
GPU Hello World with a Global

define - works on CPU, GPU

global var - only CPU. --device-- keyword.

Separate memories.

carrier used for I/O



A variable in CPU memory cannot be accessed directly in a GPU kernel.

Programmer needs to maintain copies of variables.

Typical CUDA program flow →

cuInit() → "Initialize API"

gpuErr() →

kernel<<1,32>>(apawn, strnt, bvar);

if thread is < array length, increment to character at

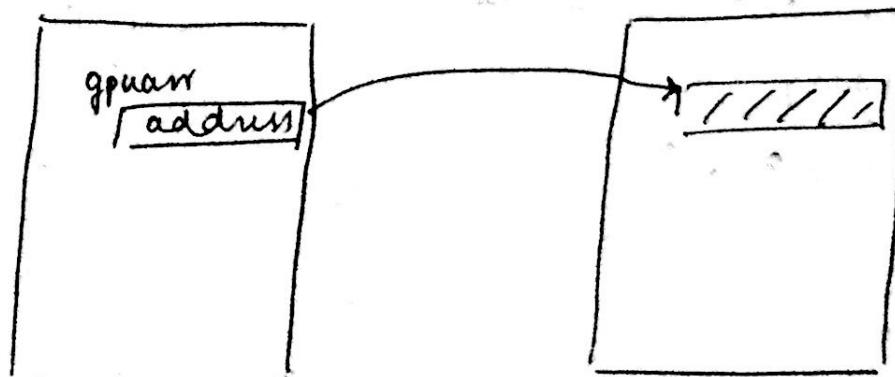
→ Hello ^{space}World [?] block

> end

for all the threads → execute the kernel but don't do anything

cpu memory

GPU memory



* gpuarr - should NOT be accessed on CPU

Same address could be for CPU or GPU.

- a) what does `cudaDeviceSynchronize()` do? Why is it not reqd? \Rightarrow
* GPU stream - queue. \rightarrow No analysis of whether GPU mem is being used

cuda Malloc	cudaMemcpy host-to-device	kernel launch	cudaMemcpy device-to-host
-------------	---------------------------	---------------	---------------------------

cuda memcpy device-to-host \rightarrow blocks -

- a) Passing stolen (cpuarr) to ~~the~~ dkernel?

\Rightarrow pass by value \checkmark okay.

Problem is only when we pass by reference.

Classwork.

1. CUDA program - initialize an array of size 32 to all zeros in parallel.

```
#include <stdio.h>
#include <cuda.h>

__global__ void dkernel(int *arr)
{
    arr[threadIdx.x] = 0;
}

int main()
{
    int cpuarr[32];
    int *gpuarr; // no need of CPU variable
    cudaMalloc(&gpuarr, 32 * sizeof(int));
    cudaMemcpy(gpuarr, cpuarr, 32 * sizeof(int), cudaMemcpyHostToDevice);
    dkernel << 1, 32 >> (gpuarr);
    cudaMemcpy(cpuarr, gpuarr, 32 * sizeof(int), cudaMemcpyDeviceToHost);
    return 0;
}
```

2. change array size to 1024 ✓

3. another kernel that adds i to array [i]

--- global add array[1] (int *arr) {

arr [iteration] += threadIdx.x; ✓

→ add in main() also --

4. change array size to 8000 --- looks okay.

[kernel kernel]

Stream - queue

No init happens, then addition..

but... dkernel <<< 1, 8000 >>>

↓
not allowed

max limit ~ 1024 or so.

dkernel <<< 8, 1000 >>>

Something like that is okay.

But code doesn't work correctly. $\text{threadIdx.x} \rightarrow$ only refers to 2nd parameter, not 1st parameter of <<<, >>>

Note: Kernel call semantics - same as C function call semantics
in terms of nested -- etc.

Thread Organization

Kernel launch : grid of threads.

Grid: 3D array of thread blocks.

(gridDim.x, blockDim.y, blockDim.z)

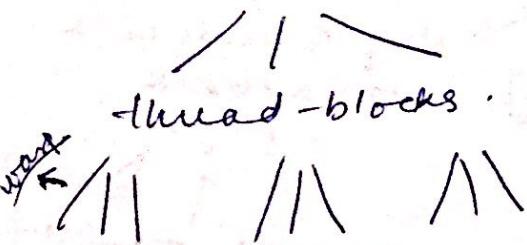
Each block \rightarrow blockIdx.x, .y, .z

A thread block is a 3D array of threads

(blockDim.x, .y, .z)

Each thread \rightarrow threadIdx.x, .y, .z

grid \Leftrightarrow kernel.



$\langle\langle 8, 1024 \rangle\rangle$
↑ ↓
threads in blocks. # threads in each block.

threads.

$\langle\langle 1, 1024 \rangle\rangle$

v/s.

$\langle\langle 2, 512 \rangle\rangle$

\rightarrow large parallelism

if it is a streaming multiprocessor, thread blocks get assigned to separate hardware \Rightarrow

20/11/2020.

One kernel - one grid. (1-1 mapping)

No performance difference for 1dim vs 2dim vs 3dim

↳ only simplifies programming ↳

- 1-5 threads get assigned to a streaming multiprocessor.
- Few 100s of 10s assigned to a thread block.
- Total few 1000 to 100 000 threads.

dim3 : structure provided by CUDA

→ qgrid dims
→ block dims

e.g. dim3 grid(2, 3, 4);
dim3 block(5, 6, 7);

dkernel << qgrid, block >>();

for 2d → dim3 qgrid(a, b, 1);
or dim2 qgrid(a, b);

and so on. . .

not 0

total # threads = $2^*3^*4^*5^*6^*7$

→ Ho of threads in a thread block
= 5^*6^*7

↑ ↑
threadIdx threadIdx

By default → 1d (single number).

threadIdx → thread ID within a block

blockIdx → block ID within a grid.

For e.g. program → Only 1 line of input.

2 3 4 5 6 7.

Q) If cond is changed to threadIdx.x == 0 ?

2D

$N \leftarrow 5, M \leftarrow 6$
 $\text{dim3 block}(N, M, 1);$
 $\text{dkernel} \ll< 1, \text{block} >> (\text{matrix});$

Though it is intended to be a matrix,
 ↘ it is passed as a single dimensional array.

↑ ↑
 1 block there many threads in the block.

How to access.

$\text{matrix}[tidx.x * \text{blockDim.y} + tidx.y]$

$a[i][j]: (i * \# \text{cols} + j)$

Q) write the kernel to initialize the matrix to unique ids.

-- global void dkernel (unsigned *matrix) {

 unsigned uid = (tidx.x * blockDim.y + tidx.y);

 matrix [uid] = uid;

3

tidx.x comes from 0, 1, 2, 3, 4

tidx.y comes from 0, 1, 2, 3, 4, 5

tidx.y	0	1	2	3	4	5
tidx.x	0	1	2	3	4	5
0	0	1	2	3	4	5
1	6	7	8	9	10	11
2	12	13	14	15	16	17
3	18	19	20	21	22	23
4	24	25	26	27	28	29

Q) can we access the matrix in a 2D way?

→ $\text{matrix}[i][j]$

Q) make it unsigned * matrix everywhere...?

→ we are not passing # columns.
 so compiler will not know how to find index...
 (M) ~~works~~

Q) unsigned * matrix [# columns].
 → works in CPU, not GPU

why `unsigned *matrix [M]` doesn't work?

HW.

→ Arrays are passed by reference.

Pay it out.

* Consecutive mem locs should hold the addresses of the M rows.

* 1 cuda malloc followed by M individual cuda mallocs for consecutive elements.
 ↳ `cudaMalloc(&a...)` ↑ `a[0]`

* dereferencing on CPU shouldn't be done.

ID.

N ← 5, M ← 6..

dkernel <<< N, M >>> (matrix);

write corresponding kernel.

④ ~~global~~ void dkernel(unsigned *matrix){
 M = matrix[0];
 int grid_id = threadIdx.x + blockIdx.x * blockDim.x;
 matrix[grid_id] = grid_id;

3

21-1-2020

takeaway: One can perform computation on a multidimensional data using a 1-d block.

`cudaMemcpy`; blocking for CPU & GPU.

`cudaMemcpy Async` - not blocking for CPU.

Launch configuration for huge data.

Good practice: ~~use~~ $\ll< - , \ge>>$

↑

Fix this param as constant, say 1024.
(Block size).

$\text{ceil}(N / \text{BLOCKSIZE})$.

↓

This itself would result in an integer division
→ need to convert to floating point division.

$\text{ceil}(\text{float}(N) / \text{BLOCKSIZE});$

could result in more no. of threads than expected.

→ add if condition in the kernel, to avoid access out of bounds.

if ($i < \text{vectorSize}$)

$\text{vector}[id] = id$

$\frac{N + BS - 1}{BS}$

- 8). Read several pairs as (x, y) words from input. For each pair of points, compute Euclidean distance in parallel. Print the maximum distance.

```
struct point { int x; float y; };
int main(int argc, char* argv[])
{
    int n = atoi(argv[1]); // # pairs of points.
    struct point *matrix; // malloc
    int i;
    for (i = 0; i < n; i++) {
        matrix[i].x = atoi(argv[2 * i]);
        matrix[i].y = atoi(argv[2 * i + 1]);
    }
}
```

-- global -- void & kernel
(int *dist, struct point *matrix);

int uid = blockIdx.x *
blockSize.x *
blockIdx.x;

dist[uid] = dist b/w points
matrix[uid].x, matrix[blockIdx.x];
matrix[blockIdx.x];

cuda_dkernel << N, N >>;
int max = 0; for loop -> update max = point(max); } }

- we could launch ~~n~~ threads. So, a kernel will take one point and find its distance from all ~~n~~ real points from it onward.
- * If we have a for loop inside the kernel, then that for loop executes sequentially.

- ~~n, n~~ threads. ~~l, NxN~~ threads
- [element ~~broadcasted~~] ~~block~~ [element ~~updated~~]. \leftarrow check.

~~blockIdx~~ = blockIdx.

~~threadIdx~~ = threadIdx.

maximum:

If we put $\text{if } (d > \max)$ inside the kernel,
~~max = d;~~ access to / write to common
race condition can occur, since ~~it~~ is not
done in an atomic manner.

Data race.

1. Shared variable / resource
2. At least 1 write operation.
3. Multiple parties / threads.
4. ?

- Max can be calculated sequentially on the CPU using for loop
- GPU: CUDA instruction: atomic maximum
 $\text{AtomicMax}(\&\max, \text{dest}) \leftarrow$ each thread (kernel).

- We can use locks also - sequential ...
- We can launch another kernel with single thread on GPU. It will calculate max sequentially. But no need to copy the vector back from GPU to CPU
- Divide & conquer mechanism.
Each thread takes a pair of distances & finds max.
keep doing this for a few more steps --
→ "reduction" technique.

24/1/2020.

Evolution of GPUs:

1995 - 2000 : Fixed function

2001 - 2005 : Programmable shaders. - variety of ways to render graphics
⇒ if it is programmable, you can use for apps other than Graphics

2006 onwards : CUDA

GPGPU : General Purpose Graphics Processing Unit.

Nvidia → Pioneers in this field.

AMD

Other vendors : Intel, Qualcomm,
ARM, Broadcom, Samsung -
not very mainstream

- * Vertices / Pixels
- * Triangulation
- * Fragments
- * Shader → Video memory.

Graphics - very good setting for parallel computations.

Interesting aspect of GPU.

Data Parallelism v/s Instruction Level Parallelism

Data Parallelism v/s Instruction Level Parallelism
eg. pixel
Each data item is being operated upon by different threads. — GPU
SCALABLE

CPU
inherent limitation due to sequential nature of code.

GPU languages.

- CUDA (compute unified device language).
 - proprietary, NVIDIA-specific
- OpenCL (open computing language)
 - universal, works across all computing devices (not only GPU)
 - CPU, FPGA etc.
- OpenACC (open accelerator)
 - universal, works across all accelerators
 - FPGA, GPU - not CPU.

2. # pragma openacc ...
 for (---) // C or C++
 _____ // implicit parallelism

compiler takes care of parallelizing C code.

II. Programmer explicitly tells what to parallelize
 → CUDA, OpenCL : explicitly parallel.

Interfaces.

Python → CUDA

Javascript → OpenCL

LLVM → PTX.

| /
 intermediate representation

Kepler family v/s Pascal family

DVFS : Dynamic Voltage Frequency Scaling.

Compute capability - versioning system
 which features are supported.

GPU RAM - much less than CPU RAM.
mainly cost issue.

top500.org - list of top supercomputers in the world.

- ranked by performance (FLOPs)
- As of Nov 2019, 1) Summit from USA (>2.4 million cores)
2) Sierra from USA (>1.5 million cores)
3) Fugaku Light from China (>10 million cores).

India - 57th position - Pratyush (Pune).

(Indian Institute of Tropical Meteorology, Pune, India).

119232 Gres

Matrix Squaring.

CPU -

void squareCPU(unsigned *matrix, unsigned *result, unsigned matrixSize) {

for (unsigned i = 0; i < matrixSize; i++) {

for (unsigned j = 0; j < matrixSize; j++) {

for (unsigned k = 0; k < matrixSize; k++) {

result[i * matrixSize + j] +=

+ matrix[i * matrixSize + k] *

matrix[k * matrixSize + j];

}

j

g

g

Q) How to parallelize?

→ each element of result matrix can be calculated in parallel
(one row, one column).

global void squareGPU (unsigned *matrix, unsigned *result, unsigned matrixSize) {

unsigned ii = blockIdx.x * 32 + threadIdx.x; ↳ seems okay... -

unsigned jj = threadIdx.y * 32 + threadIdx.y;

for (unsigned kk = 0; kk < 32; kk++) { ↳ finds kk * 32

24-1-2020

Matrix squaring - Parallelize i loop (outer loop)

Matrix multiplication

for (int i = 0; i < matrixSize; i++) {
 for (int j = 0; j < matrixSize; j++) {
 result[i][j] = 0;
 for (int k = 0; k < matrixSize; k++) {
 result[i][j] += matrix[i][k] * matrix[k][j];
 }
 }

for (int i = 0; i < matrixSize; i++) {
 for (int j = 0; j < matrixSize; j++) {
 result[i][j] = 0;

for (int k = 0; k < matrixSize; k++) {
 result[i][j] += matrix[i][k] * matrix[k][j];
}

result[i][j] = result[i][j];
t = matrix[0:matrixSize-1:k]
* matrix[k:matrixSize-1:i]

}

2

3

→ Slower than CPU

⇒ Parallelization of i loop alone is not able to compensate for the GPU's lower clock frequency.

time about

source of contention for memory

start time + n * clock

kernel call launch

cuda device switch

end time + n * clock

printme (" ", start, end)

system.h

address of data

St = clock()

Kernel launch

cuda device switch

end = clock()

printme (" ", end - St, 0)

GPU - lower clock frequency than CPU's.

So sequential code on GPU is slower than on CPU

Parallelize both i and j loops

```

dkeras<-N,N>> { matrix, result, N },
  -> op2 <-> void dkeras (matrix, result, matrixSize) {
    unsigned id = blockIdx.x * blockDim.x + threadIdx.x;
    unsigned ii = id / matrixSize;
    unsigned jj = id % matrixSize;
    } } HW what happens when
    } you interchange ii,jj?
    for (unsigned kk=0; kk<matrixSize; ++kk) {
      result [ii*matrixSize+jj] += matrix [ii*matrixSize+kk]
                                * matrix [kk*matrixSize+jj];
}

```

3

Sometimes we deliberately make a thread more sequential to improve performance.

When we try to parallelize k loop also - race condition may arise. Threads will have to coordinate among each other. \Rightarrow Overhead will result.

ii,jj loops are "embarrassingly parallel" - scales very well
kk loop is not

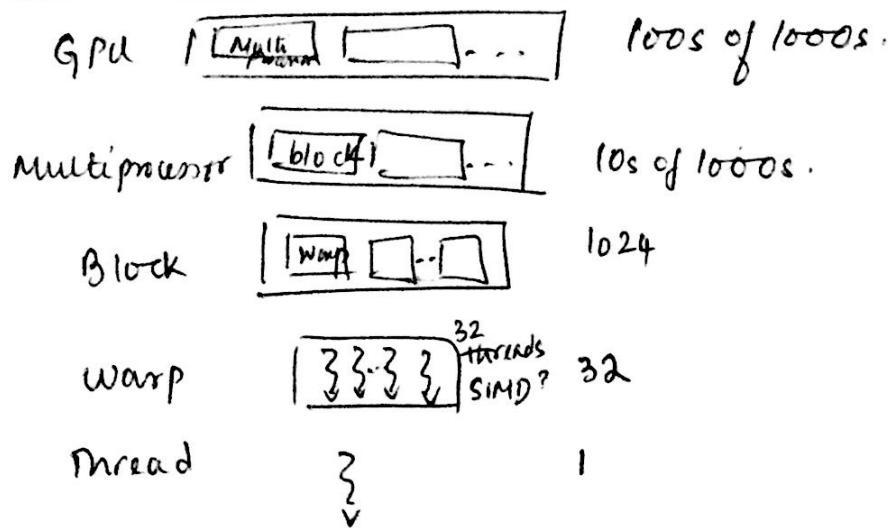
When no. of cores increased beyond a limit, Then physical architecture/organization comes into picture.

\rightarrow Systemic aspects come into picture.

PRAM model - theoretical model.

von Neumann architecture.

GPU computation hierarchy.



Each core in a multiprocessor executes one thread.

Multiple grids could be simultaneously executing on a GPU

27/11/2020.

Warp : SIMD. - Single Instruction Multiple Data

- * Warp threads are fully synchronized.
- * Implicit barrier after each step/instruction.

Lock step execution is only true for threads within a warp.

- * Hardware takes care of adding the NOPs.

* Thread privatisation
Modify local copies in each thread; after barrier, merge the computations of each thread.

When diff warp threads execute different instructions, threads are said to diverge. Hardware executes threads satisfying some condition together, ensuring that other

threads execute no-op. This adds sequentiality to the execution.
⇒ "thread divergence"

How to avoid?

1) Carefully remove if conditions

2) Reorganize data / computation. - separate chunks.

* simple if conditions are not as much ~~as~~ a problem as
for loops, with different # iterations for diff warptracks

→ Reduce disparity across amt of work that each warp
thread is doing. (bucket sort, in some sense).

→ reduces thread divergence within a warp.

* static * loadtime * dynamic
bucket sort inspect data,
 then do next
 part... can't use
 any sorting
 method..

degree of Divergence.

DOD for a warp is the # of steps reqd to complete 1 instruction
for each thread in a warp.

without any thread divergence, $DOD = 1$

for fully divergent code, $DOD = 32$.

CW: Write a code to achieve $DOD = 4$.

int* matrix, int N

```
-- global-- void __kernel void L2 {
    int uid = blockIdx.x * blockDim.x + threadIdx.x;
    if (uid % 4 == 0) {
        matrix[uid % N] = uid;
    } else if (uid % 4 == 1) {
        matrix[uid % N] = uid + 1;
    } else if (uid % 4 == 2) {
        matrix[uid % N] = uid + 2;
    } else if (uid % 4 == 3) {
        matrix[uid % N] = uid + 3;
    }
}
```

Q) Calculate degree of divergence of :

case 0

$$\text{vector}[id] = 0$$

case 1

$$\text{vector}[id] = \text{vector}[id]$$

case 2

$$\text{vector}[id] = \text{vector}[id-2]$$

case 3

$$\text{vector}[id] = \text{vector}[id+3]$$

case 4

$$\text{vector}[id] \neq 4 + 4 + \text{vector}[id].$$

case 5

case 6.

$$\text{vector}[id] = 5 - \text{vector}[id]. \quad \text{vector}[id] = \text{vector}[6].$$

case 7

$$\text{vector}[id] = 7 + 7.$$

case 8

$$\text{vector}[id] = \text{vector}[id] + 8$$

case 9

$$\text{vector}[id] = \text{vector}[id]^9;$$

default.

Answer - dependence
Compiler Optimisations

11 10 9

gcc by default -O3

- O1 } for performance
- O2 } optimization.
- O3 }

nvcc

→ by default -O3

① There can be conditions in the kernel. But there needn't be divergence.

if (vectorsize < N) S1; else S2;

② Conditions evaluating to diff. truth values.

→ Need not lead to divergence

e.g. if (id / 32) S1; else S2;

TL no divergence within a warp.
∴ evaluates to same truth value

③ Conditions evaluating to the same truth values within a warp.

28-01-2020

Memory.

CUDA Memory model overview.

- Within a thread block : Shared memory.
 - Blw blocks - share global memory.
- Shared memory - similar to L1 cache.
L1 or scratchpad cache.

We can decide, to a certain extent, what goes in this 'cache'. ^{Registers} ^{16-32 KB?} ^{on chip}
Multiprocessor has a separate register file. ^{separate from core-specific registers}
Registers are also available ^{for} each thread. ^{split among the threads}
~~can be~~

registers is good for a particular thread if we give the "right" no. of threads -
else register spills ^{goes into global memory} → poor performance.

Q) How to check if # registers is good enough?

NVCC gives a warning...

Global / video memory

→ main means of communicating data blw host & device
(PCI express bus).

→ contents visible to all GPU threads.

→ longer latency. (400-800 cycles)

→ Throughput ~200 GBPS.

Texture Memory. (smaller)

→ read only (12 KB)

→ ~ 800 Gbps

→ Optimized for 2D spatial locality

Bucket Sort.

We have marks from 0-100.
Want to divide into buckets
0-10, 10-20 ... 90-100.

- 1) Sort in $O(n \log n)$, then do a linear pass, $O(n)$.
- 2) Push each element/mark into a separate vector
(10 diff. vectors: 0-10, 10-20 ... 90-100).
→ $O(n)$ time.

* cudaGetLastError() to catch errors on GPU.

Constant Memory.

→ read only
(64 KB) -

SM: Streaming Multiprocessor

- L2 Cache.
 - 768^{KB} cache
 - Shared among SMs
 - Fast atomic instructions can be supported.

- L1 / shared memory.

- Configurable 64 kB per SM.
 - 16 kB shared + 48 kB L1 or vice versa
 - low latency (20 - 30 cycles)
 - high bandwidth ($\sim 1\text{ TBPS}$)
- Registers. unified,
 - 32 K registers per SM.
 - Max ~ 21 regs per thread.
 - very high bandwidth ($\sim 8\text{ TBPS}$)

Bandwidth.

- Wide data bus rather than fast data bus
- Parallel data transfer.
- Techniques to improve bandwidth:
 - Share/reuse data
 - + Data compression. (decompress locally)
 - Recompute rather than store + fetch

Latency

- Time reqd for I/O / data transmission.
- Ideally, it should be zero. But not possible. So it should be minimized as much as possible.
 - Processor should have data available in no time.
 - memory I/O is bottleneck.
- Latency can be reduced using cache.
 - CPU: Regs, L1, L2, L3, L4, RAM
 - GPU: small L1, L2, many threads
- Latency hiding on GPUs is done by exploiting massive multithreading, rather than ILP.

Individual latencies can be large. But overall throughput is maximized. The GPU keeps itself busy all the time - as far as possible.

hiding latency via pipelining

* May ~~not~~ lead to super-linear speedup.

CPU - ILP : take another instruction
GPU - take another thread block.

Locality

- Only a small cache is available per thread.
- In GPU, ^{another form of} spatial locality is critical.

\Rightarrow Across threads

All threads in a thread block access their L1 cache.

Pascal GPU

64 KB L1 cache
16 KB + 48 KB /
48 KB + 16 KB ·
(4 MB scratchpad)

29-1-2020

Locality.

1. Spatial

If $a[i]$ is accessed, $a[i+k]$ would also be accessed soon.

```
for (i=0, i<N; i++)
```

```
    a[i] = 0;
```

2. Temporal

If $a[i]$ is accessed now, it would be accessed again soon.

```
for (i=0, i<N; i++) {  
    a[i] = i;  
    a[i] += N;  
    b[i] = a[i] * a[i];  
}
```

$a[i]$ - temporal locality

$b[i]$ - no temporal locality

These localities are applicable to both CPUs & GPUs, but more so on CPUs.

(Cw)

```
for (i=0, i<M; ++i)
```

I for (j=0, j<N; ++j)

```
    for (k=0, k<P; ++k)
```

```
        C[i][j] += A[i][k] * B[k][j];
```

II

```
for (i=0, i<M; ++i)
```

```
    for (k=0, k<P; ++k)
```

```
        for (j=0, j<N; ++j)
```

```
            C[i][j] += A[i][k] * B[k][j];
```

On CPU, II takes half as much time as I.

$i \rightarrow k \rightarrow j$

Access to k should be after i . ✓ in both cases

But j should be after k → [poly in 2nd one]

So 2nd way is more efficient ✓

DAG-topological order of access patterns.

Much of dependences : construct DAG ; use topo order for efficient cache accesses.

Q) what about GPU?

Memory coalescing.

to avoid warp stalling
(recall: barrier--)

If consecutive threads access words from the same block of 32 words, their memory requests are clubbed into one.

e.g. If $a[id]$ is accessed by thread 'id', then $a[id+1]$ is also brought into the cache of thread 'id + 1'.

"Memory requests are coalesced".

→ "stride": $a[i] \rightarrow a[20 * i + 1]$. Strided access pattern.

as opposed to $a[i+1]$

- Memory coalescing can be effectively achieved for regular programs (such as dense matrix operations).

$a[id] \checkmark \quad a[id+5] \checkmark$ (no need to be at boundary of 32)

- CUDA can be associated with C as well as Fortran

↑ ↑
row major access column major access.

- Coalesced - doesn't matter which order the memory is accessed

- Not only read accesses, but also write accesses are coalesced.

- Vectorization on Intel processors \rightarrow somewhat similar to this

* Large instruction / Vector instruction

\rightarrow same instruction, multiple data

\rightarrow sequential program

\rightarrow handled by hardware.

Intel compiler: `icc` \rightarrow automatically vectorizes "good" access patterns.

- Degree of coalescing:

- DOC is the inverse of (32 -) the no. of memory trans reqd for a warp to execute an instruction.

31

e.g. $a[\text{ThreadIdx}, x]$ has DOC of ~~32 (31?)~~

$a[\text{rand}()]$ likely has a DOC of 0.

CW: Write a kernel to vary the degree of coalescing from 0 to 32 based on an input argument.

```
--global-- void kernel (int doc) {  
    1) access a [ThreadIdx.x * doc]  
    2) declare a [doc];
```

3

31-1-2020 Tutorial Session

03-02-2020.

Memory Coalescing.

- Gallored - deterministic sequential access
- Strided access - sequential access but in chunks.
 $a[\text{rand}()]$
- random access - $a[\text{input}[id]]$: input-dependent pattern.
↳ e.g. bucket sort -
graphs representation.
edges/ neighbours of a vertex.

① CPU - \rightarrow spatial locality
- Each thread should access consecutive elements of a chunk (strided).
- Array of Structures (AoS). has better locality.

② GPU - \rightarrow coalescing

- A chunk should be accessed by consecutive threads (coalesced).
- Structure of Arrays (SoA) has a better performance.

AoS v/s SoA.

```

AoS
struct node {
    int a;
    double b;
    char c,
};

struct node allnodes[N];

```

Expectation: When a thread accesses an attribute of a node, it also accesses other attributes of the same node.

LOCALITY
BETTER ON CPU

SoA.

```

struct node {
    int alla[N];
    double allb[N];
    char allc[N];
};


```

?

Expectation: When a thread accesses an attribute of a node, its neighbouring thread accesses the same attribute of the next node.

COALESCING.

BETTER ON GPU.

Q) Write codes for the two types using cudaMemcpy.

Note that all arrays would be pointers.

AoS.

```

struct node {
    int a;
    double b;
    char c;
};


```

?

struct node* allnodes;

struct node* datanodes;

```

cudaMalloc(&datanodes,
N*(sizeof(
struct node)));

```

cudaMemcpy(datanodes, allnodes,

N*(sizeof(struct node)),

cudaMemcpy(datanodes, allnodes,

SoA. ??

struct node * allnode, * datanode;

cudaMalloc(&datanode,

↓

int*

double*

char*

int

double

char

not aligned

cudaMalloc(&datanode, size of int * N);

cudaMalloc(&datanode, size of double * N);

cudaMalloc(&datanode, size of char * N);

cudaMemcpy(datanode, allnode, size of int * N, cudaMemcpyHostToDevice);

cudaMemcpy(datanode, allnode, size of double * N, cudaMemcpyHostToDevice);

cudaMemcpy(datanode, allnode, size of char * N, cudaMemcpyHostToDevice);

and so on

② ~~for i = 0 to k
 a[i] = 0.0f;
 printf("%f", a[i]);
 scanf("%f", &a[i]);
 cout << a[i];~~
 cout << endl;

③ struct node
 { int a[3][3];
 double b[3][3];
 float c[3][3];
 node *next;

node *head = new struct node();

cout << "size of (size of (struct node))";
 cout << endl;

cout << "size of (size of (int) * N)";
 cout << endl;

④ cout << "size of (size of (struct node));
 cout << endl;

cout << "size of (size of (double));
 cout << endl;

cout << "size of (size of (float));
 cout << endl;

cout << "size of (size of (int));
 cout << endl;

cout << "size of (size of (double));
 cout << endl;

* malloc() is available on GPU.

✓
SOA takes almost $\frac{1}{3}$ rd the time as AOS.

24/02/2020

classwork

copy a linked list from CPU to GPU

- Each node contains roll no., name, facad.

In main program
ptr = malloc(sizeof(node))

cuda malloc & free

node malloc & free

ptr is on CPU

its value is an address on GPU

blocks

Start on GPU, ?

ptr on CPU which points to GPU memory

cuda malloc & free

CW.

head

→ start node



ptr, head, gpu ptr, gpu head, a
while(ptr != NULL & gpu ptr != NULL) {

struct node {

char * rollNo;

char * name,

char * facad,

struct node * next;

cudaMemcpy(gpu ptr, ptr, sizeof(struct node));

cudaMalloc(&gpu ptr, sizeof(struct node));

→ name

→ facad

→ next

3

cudaMemcpy(gpu head, head, sizeof(struct node));

free(ptr); free(head);

struct node * head, *ptr;

struct node * gpu head, *gpu ptr;

→ name

→ facad

→ next

M

Recursive code :

final head of
returns parallel to linked list
in GPU

```
node * getLinked (node *root) {
    if (root == NULL) return NULL;
    node * ptr; on CPU, points to a mem in GPU
    cudaMalloc (&ptr, size of node * nodes);
    node * n_ptr = getLinked (root -> next);
    struct node * copy = root;
    copy -> next = n_ptr; copy data in root along with ptrs.
    cudaMemcpy (ptr, copy, ... );
    return ptr;
}
```

3

HW: Strlen parallel implementation. — without array out of bounds exception

Idea: Break up the given string into some number of small chunks. Each chunk's length is found sequentially by a thread. If '0' is found we know that the entire string is scanned. One thread