# Pinned memory
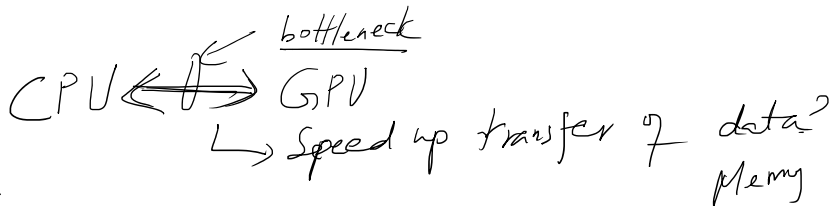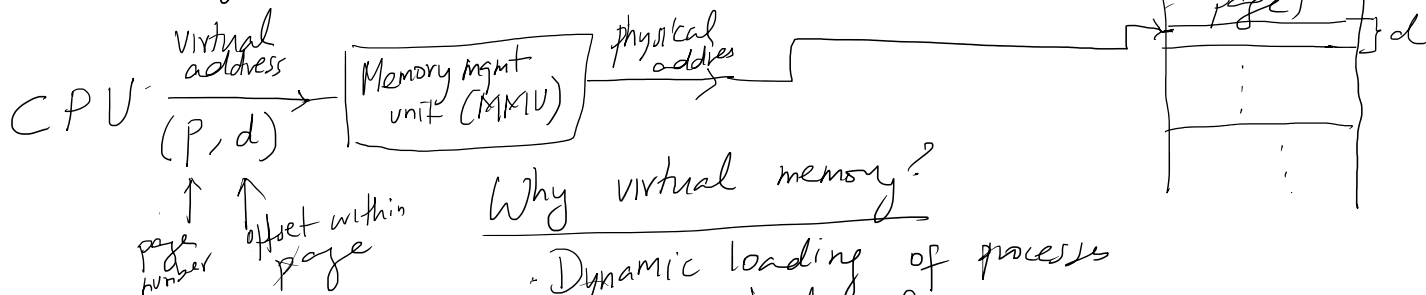
- Pinned or page-locked memory improves achieved bandwidth between host and device
  - Enables zero-copy transfers
- Use *cudaHostAlloc()* to request pinned pages from the Operating System
- Code example: *page_locked_memory*
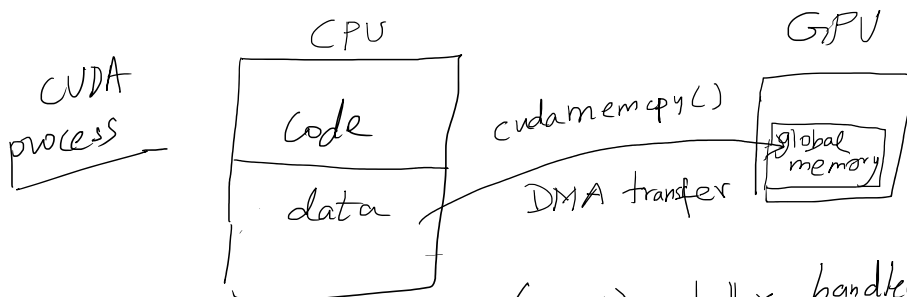
, Page locked memory / Pinned memory

CPU ⟷ GPU   bottleneck
  ↳ Speed up transfer of data?
  Memory

4KB { page 0 (P) }   address
    page 1        ∟ d

· Virtual memory ·

CPU  $\xrightarrow{\text{virtual address}}$ | Memory mgmt unit (MMU) | $\xrightarrow{\text{physical address}}$

$(P, d)$

↑ page number    ↑ offset within page

Why virtual memory?
· Dynamic loading of processes
  ↳ demand paging
Address space protection

Overhead due to address translation

CPU                    GPU
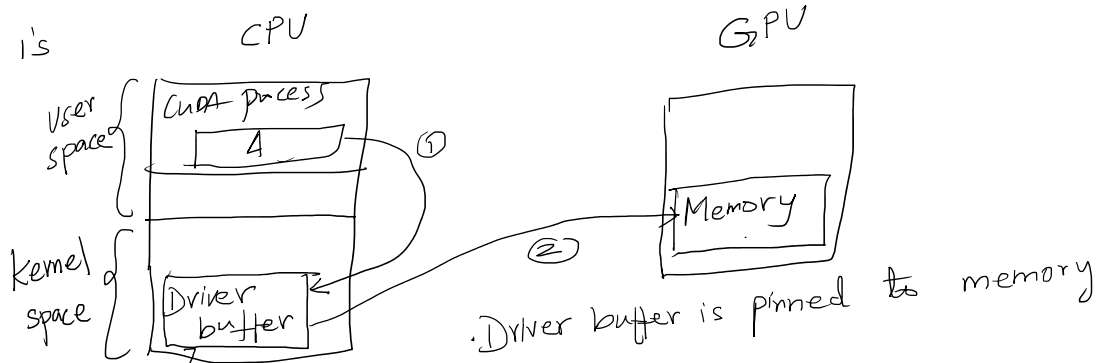
CUDA process ___  | Code |   cudamemcpy()  | global memory |
                  | data |  DMA transfer

Direct memory access (DMA) controller handles data transfer
· src   address (host)
· dest  address (device)    ‖ → DMA happens in background
· # of bytes to transfer      → Frees up CPU for other jobs.

· DMA transfer is a two step process

CPU                              GPU

user space { | CUDA process | 4 |   ① |
kernel space { | Driver buffer |      | Memory |   ②

· Driver buffer is pinned to memory

# Zero-copy transfer:

Use cudaHostAlloc() to request pinned memory from operating system

obtained using cudaHostAlloc()
$\Downarrow$
pinned to physical memory

CUDA process
Buffer
Kernel

CPU

GPU
Memory

Upto 2x speedup of CPU $\longleftrightarrow$ GPU data transfer

# CUDA streams

Code example: *streams*

So far: data is transferred in full
before computation begins

$$CPU \longrightarrow GPU$$
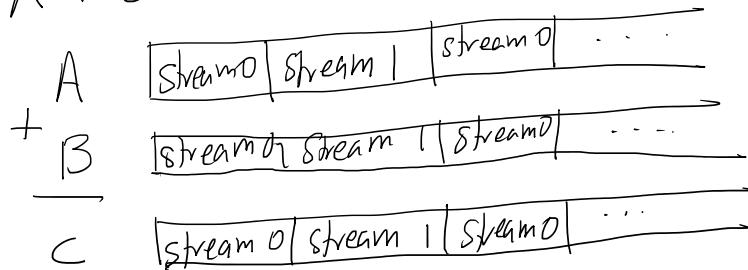
Streaming model:

1) Transfer data in smaller chunks
2) GPU can start processing chunk
3) Continue to transfer additional chunks to GPU
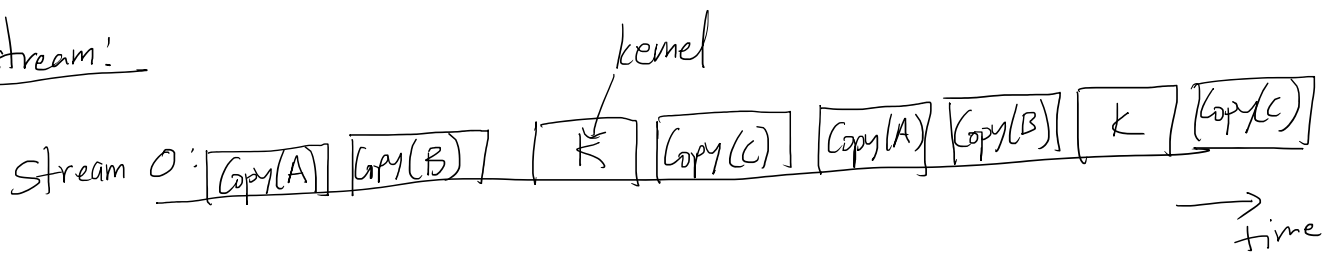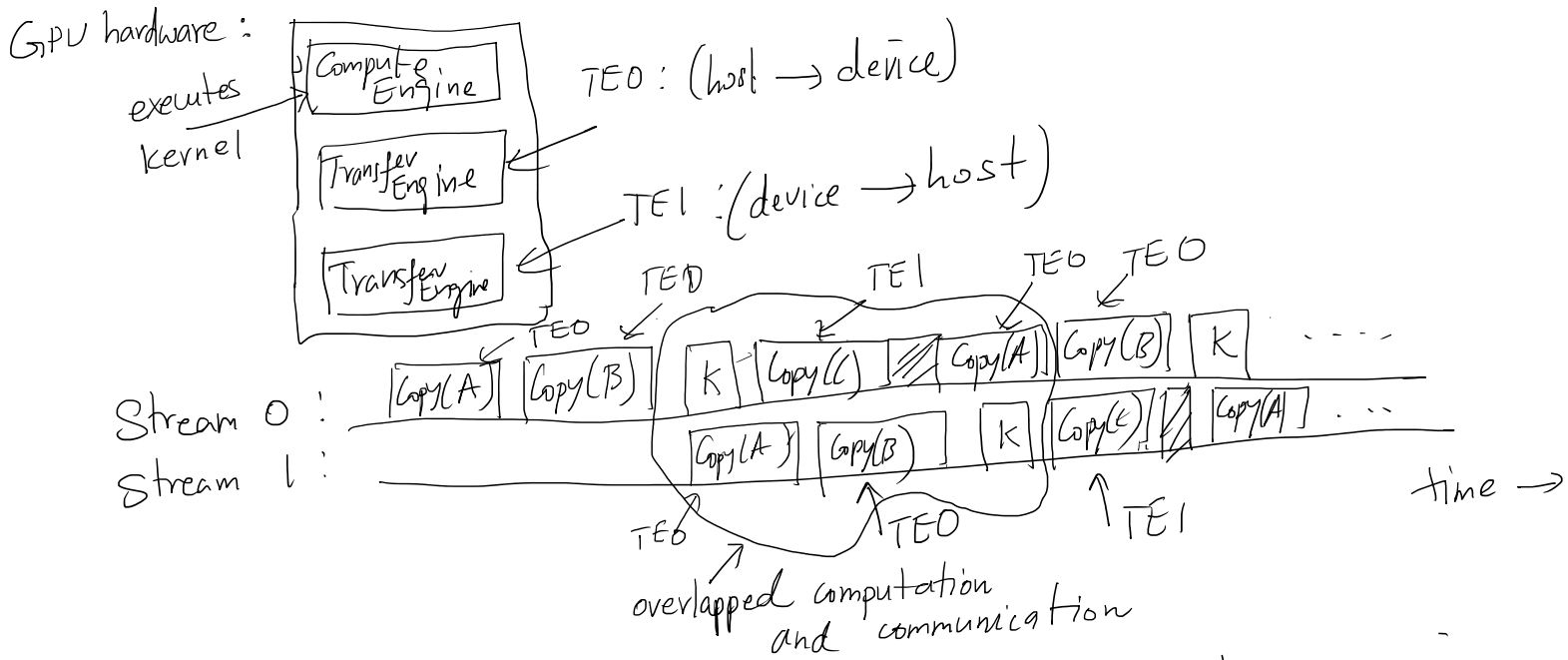4) Overlap computation + communication

Streaming source $\longrightarrow$ GPU $\xrightarrow{\text{results}}$ CPU
(sensor)
analysis ↑     ↖ Operate GPU in
                   streaming mode

Vector addition:

$$C = A + B$$

$$\begin{array}{c} A \\ + B \\ \hline C \end{array}$$

| Stream 0 | Stream 1 | Stream 0 | . . . |

| Stream 0 | Stream 1 | Stream 0 | . . . |

| Stream 0 | Stream 1 | Stream 0 | . . . |

Single stream:

kernel

Stream 0: | Copy(A) | Copy(B) | K | Copy(C) | Copy(A) | Copy(B) | K | Copy(C) |

$\longrightarrow$ time

GPU hardware:

executes kernel

Compute Engine

Transfer Engine

Transfer Engine

TE0: (host → device)

TE1: (device → host)

Stream 0:
Stream 1:

Copy(A) | Copy(B) | K | Copy(C) | ▨ | Copy(A) | Copy(B) | K · · · ·

Copy(A) | Copy(B) | K | Copy(C) | ▨ | Copy(A) · · ·

TE0   TE1   TE0   TE0

TE0   TE0   TE1

time →

overlapped computation and communication

· All operations are done in async fashion using pinned memory.

# Unified memory

https://devblogs.nvidia.com/unified-memory-cuda-beginners/

https://devblogs.nvidia.com/maximizing-unified-memory-performance-cuda/

Code example: *unified_mem*

Data transfer between the memories is managed by the CUDA runtime

CPU    GPU
Memory    Memory

Page migration between CPU and GPU:

16 kB    CPU

A : {
page 0
page 1
page 2
page 3
}

Page table
valid dirty

| | valid | dirty |
|---|---|---|
| page 0 | 1 | 0 |
| page 1 | 1 | 0 |
| page 2 | 1→0 | 1→0 |
| page 3 | 1 | 0 |

Launch kernel →

GPU
V  D

| | V | D |
|---|---|---|
| page 0 | 0 | 0 |
| page 1 | 0 | 0 |
| page 2 | 0→1 | 0 |
| page 3 | 0 | 0 |

Suppose kernel accesses data in page 2:

1) Page fault

2) Page 2 migrated from CPU to GPU

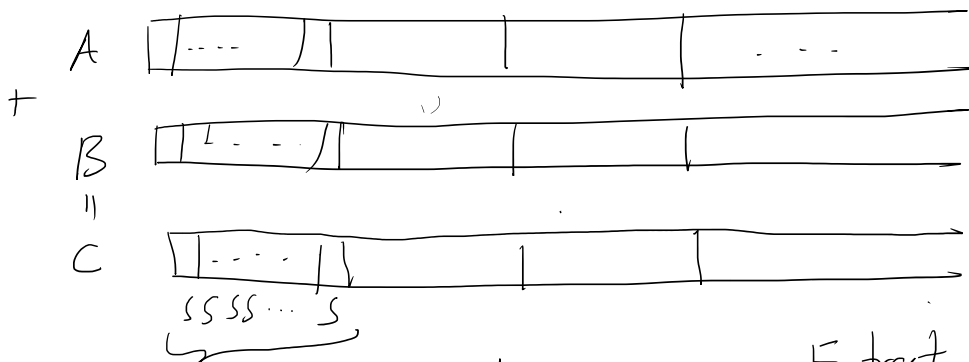3) Page tables updated on GPU and CPU

# Using CUDA with OpenMP

Code example: *cuda_omp*

OpenMP threads on CPU

$C = A + B$

- Extract coarse-grained parallelism using OpenMP

$S_0$    $S_1$    $S_2$

$A$

$+$

$B$

$=$

$C$

$SSSS \cdots S$
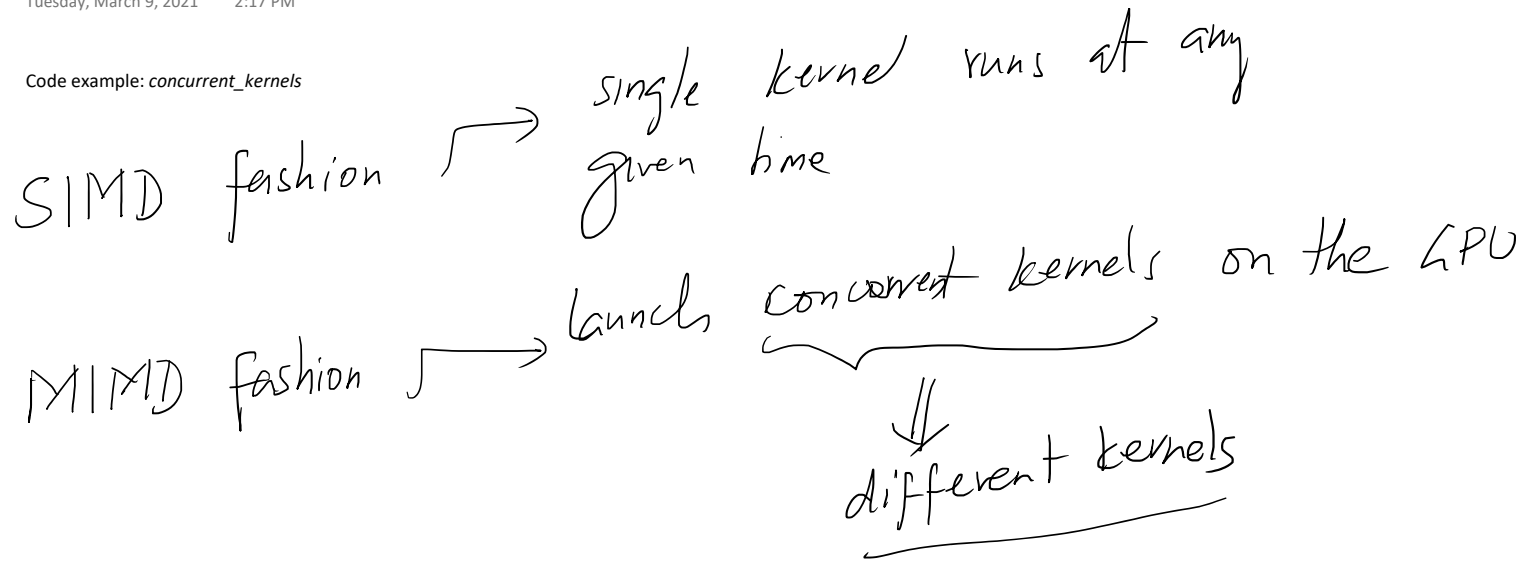
↑ CUDA threads per openMP thread

. Extract fine-grained parallelism using CUDA

# Concurrent kernel execution

Code example: *concurrent_kernels*

SIMD fashion → single kernel runs at any
given time

MIMD fashion → launch concurrent kernels on the GPU

⇓
different kernels

# CUDA BLAS examples

https://docs.nvidia.com/cuda/cublas/index.html

Single precision AX plus Y (SAXPY)

Vector dot product (dot)

Single precision general matrix-matrix multiplication (sgemm)

Single precision general matrix-vector multiplication (sgemv)

See CUBLAS documentation on BBLearn

BLAS : Basic linear algebra subroutines

# Vector processing on the CPU

Thursday, March 11, 2021     8:43 AM

- *Streaming SIMD Extensions* (*SSE*) is a single instruction, multiple data (SIMD) instruction set extension to the x86 architectureop
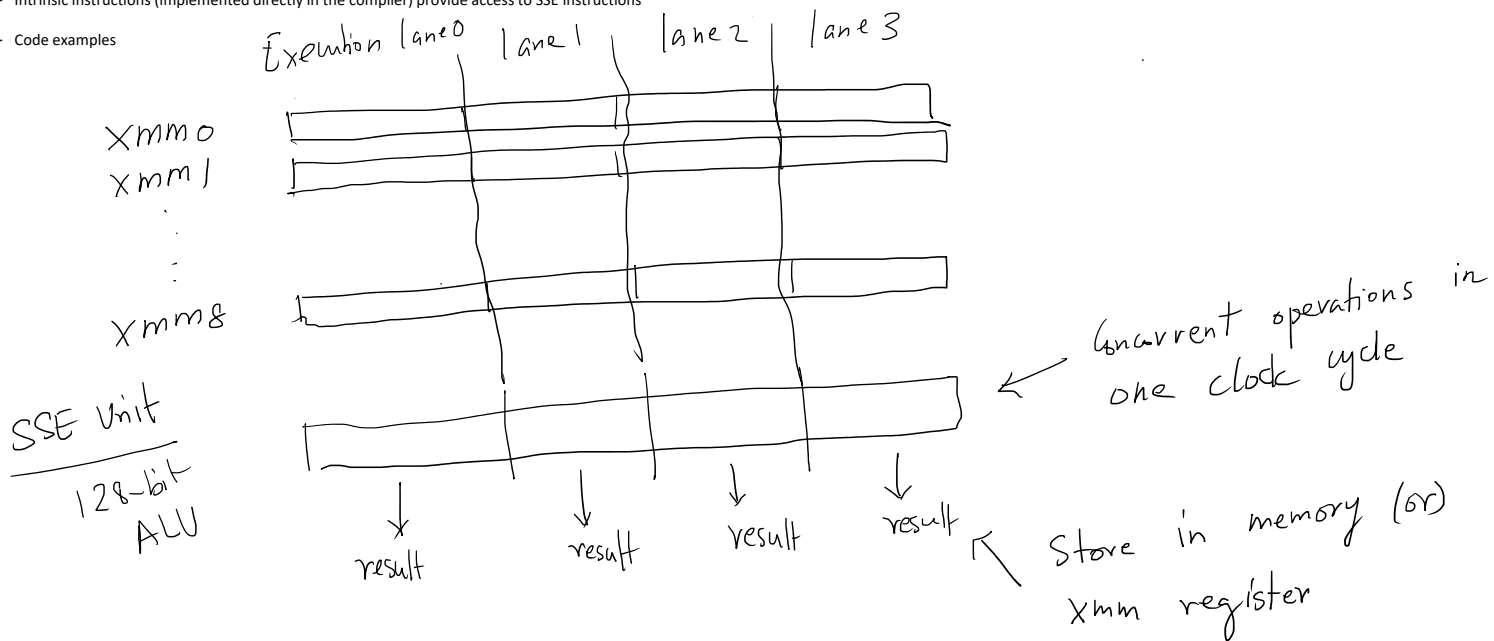
- SSE2:
  - Eight new 128-bit registers known as XMM0 through XMM7
  - XMM registers can be configured to hold
    - Four 32-bit single-precision floating-point numbers or
    - Two 64-bit double-precision floating-point numbers or
    - Two 64-bit integers or
    - Four 32-bit integers or
    - Eight 16-bit short integers or
    - Sixteen 8-bit bytes or characters
  - The ISA supports both scalar and packed scalar (vector) instructions
    - Memory-to-register/register-to-memory/register-to-register data movement
    - Arithmetic operations (add. Subtract, multiply, divide)
    - Bit-wise logical operations
    - Compare and shuffle
    - …

- More recent developments: *Advanced Vector Extensions* (*AVX*)
  - AVX uses sixteen YMM registers, each of width 256 bits, to perform SIMD operations
    - Eight 32-bit single-precision floating point numbers or
    - Four 64-bit double-precision floating point numbers
  - AVX-512 uses 32 512-bit registers (ZMM0-ZMM31) for SIMD operations

- Intrinsic instructions (implemented directly in the compiler) provide access to SSE instructions

- Code examples

$$A \cdot x = b \implies U x = b'$$

$n \times n$   $n \times 1$   $n \times 1$

$U$ upper triangular matrix

$b'$

$$\begin{bmatrix} 1 & 1 & \cdots & \\ & & \ddots & \\ \text{zeros} & & & a \end{bmatrix} \begin{bmatrix} x \\ x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{n-2} \\ b_{n-1} \end{bmatrix}$$

## Gaussian elimination

$$3x + 5y + 2z = 19$$
$$2x + 3y + z = 11$$
$$x + 2y + 2z = 11$$

$$A \begin{bmatrix} 3 & 5 & 2 \\ 2 & 3 & 1 \\ 1 & 2 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 19 \\ 11 \\ 11 \end{bmatrix} \begin{matrix} R1 \\ R2 \\ R3 \end{matrix}$$ $b$

### Back substitution

$$x_{n-1} = b_{n-1}$$
$$x_{n-2} + a \cdot x_{n-1} = b_{n-2}$$
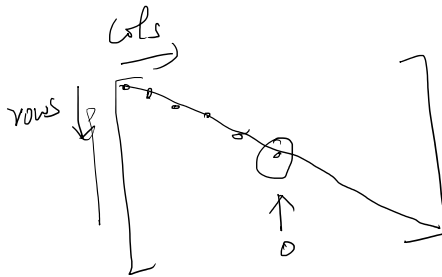$$x_{n-2} = b_{n-2} - a \cdot x_{n-1}$$
$$\vdots$$

1) Division step: $R1/3$

$$A' = \begin{bmatrix} 1 & 5/3 & 2/3 \\ 2 & 3 & 1 \\ 1 & 2 & 2 \end{bmatrix} \qquad b' = \begin{bmatrix} 19/3 \\ 11 \\ 11 \end{bmatrix}$$

2) Elimination step

$R2 = R2 - 2R1$   $\begin{bmatrix} 1 & 5/3 & 2/3 \\ 0 & -1/3 & -1/3 \\ 0 & 1/3 & 4/3 \end{bmatrix}$   $b' = \begin{bmatrix} 19/3 \\ -5/3 \\ 14/3 \end{bmatrix} \begin{matrix} R1 \\ R2 \\ R3 \end{matrix}$

$R3 = R3 - R1$

3) Division step: $R2 / -1/3$

$$A' = \begin{bmatrix} 1 & 5/3 & 2/3 \\ 0 & 1 & 1 \\ 0 & 1/3 & 4/3 \end{bmatrix} \qquad b' = \begin{bmatrix} 19/3 \\ 5 \\ 14/3 \end{bmatrix}$$

4) Elimination step:

$R3 = R3 - \frac{1}{3}R_2$   $A' = \begin{bmatrix} 1 & 5/3 & 2/3 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$   $b' = \begin{bmatrix} 19/5 \\ 5 \\ 3 \end{bmatrix}$

5) Division:   $U = \begin{bmatrix} 1 & 5/3 & 2/3 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$   $b' = \begin{bmatrix} 19/5 \\ 5 \\ 3 \end{bmatrix}$

Back substitution :

$$z = 3$$
$$y + z = 5$$
$$y = 5 - 3 = 2$$
$$x + 5/3 y + 2/3 z = 19/5$$
$$x = 1$$

- Numerical stability
  - Not stable

- Exceptions on GPU
  - Division by zero

cols

rows

- Silent exception

$$a/o \Rightarrow NaN$$
$$NaN + NaN = \overline{NaN}$$

---

Host-side code :

```
Transfer A to device
for (k=0; k<n ; k++) {     // Process rows
    Call division kernel (A, k, n)
        sync.
    Call elimination kernel (A, k, n)
}       sync.
Read A from device.
```
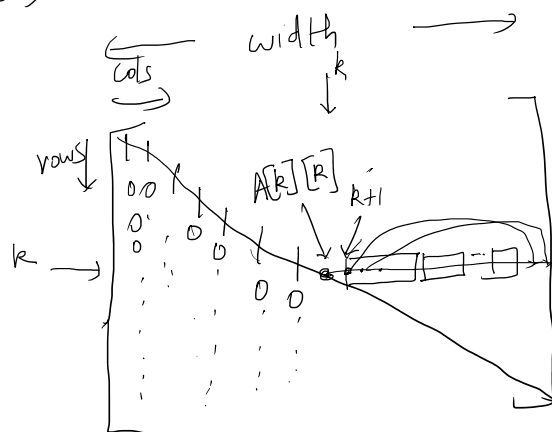
Device code:
- Design considerations:
  1) Fewer number of thread blocks
  2) Coalesced access to global memory

Division:

width
k
cols
rows

k →

A[k][R]   k+1

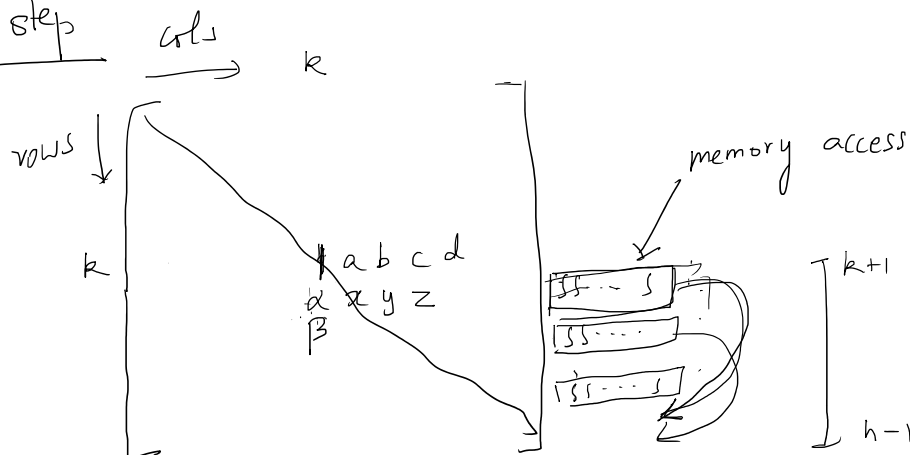#thread blocks: m
thread block size: (p,1,1)
grid (m, 1)

Kernel code:

$tid = blockIdx.x * blockDim.x + threadIdx.x$

$stride = gridDim.x * blockDim.x$

$pivot = A[k*width+k]$

while (tid < n) {
  $elem = A[k*width+k+1+tid]$
  $elem = elem / pivot$
  $tid = tid + stride$
}

if (blockIdx.x == 0 &&
   threadIdx.x == 0)
   $A[k*width+k] = 1$

Elimination step

cols
k

rows ↓

k

a b c d
α x y z
β

memory accesses are coalesced

k+1

h-1

# kernel :

$row = k + 1 + blockIdx.x$

```
while ( row < n ) {              // stride accross rows
    tid = threadIdx.x
    while ( tid < n ) {          // stride along row
        x = A[row * width + k+1 + tid]
        x = x - α·a ← A[k * width + k+1 + tid]
                  ↑ A[row * width + k]
        tid = tid + blockDim.x
    }
    sync
    if (threadIdx.x == 0)
        α = 0
    row = row + gridDim.x
}
```