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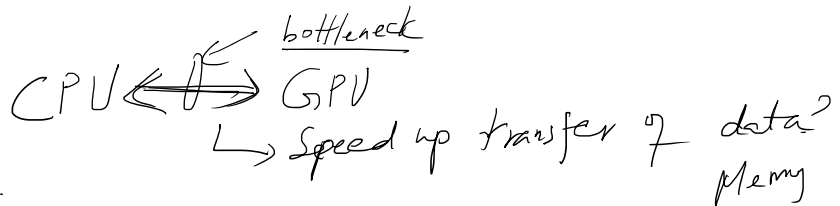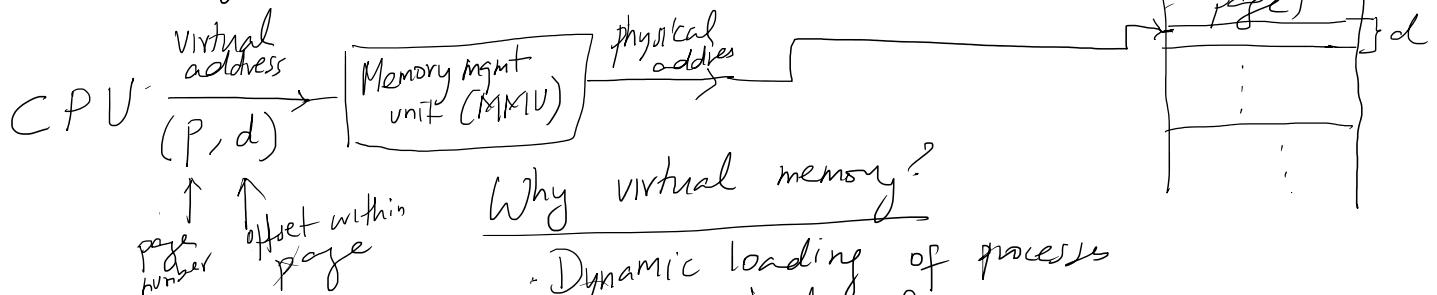
- Pinned or page-locked memory improves achieved bandwidth between host and device
  - Enables zero-copy transfers
- Use *cudaHostAlloc()* to request pinned pages from the Operating System
- Code example: *page_locked_memory*
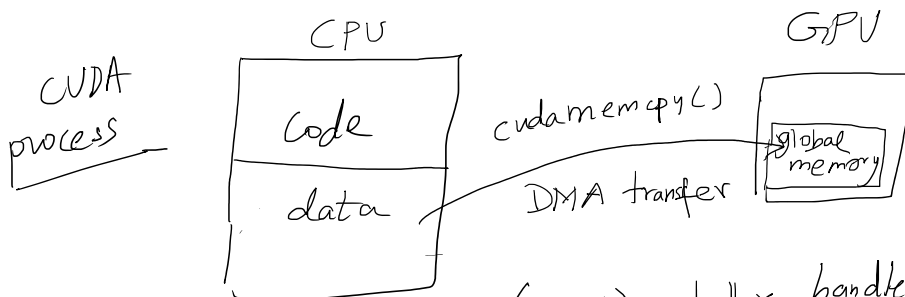
, Page locked memory / Pinned memory .

CPU $\longleftrightarrow$ GPU   ← bottleneck

$\rightarrow$ Speed up transfer of data?

· Virtual memory .

CPU $\xrightarrow[\text{(P, d)}]{\text{virtual address}}$ | Memory mgmt unit (MMU) | $\xrightarrow{\text{physical address}}$

Memory   address

4KB { page 0 (P) }   address

page 1 ← d

↑ page number   ↑ offset within page

Why virtual memory?

· Dynamic loading of processes
  └ demand paging
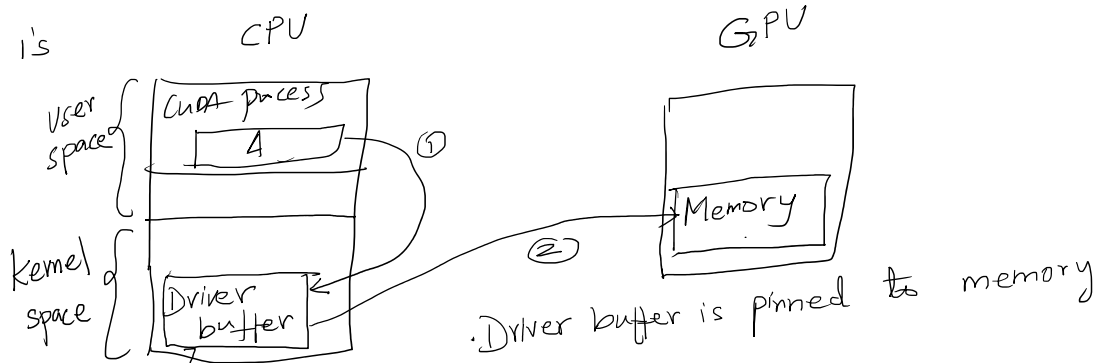  Address space protection

Overhead due to address translation

CUDA process ___

CPU

| Code |
| data |

cudamemcpy()   $\rightarrow$ global memory   GPU

DMA transfer

Direct memory access (DMA) controller handles data transfer
· src    address (host)
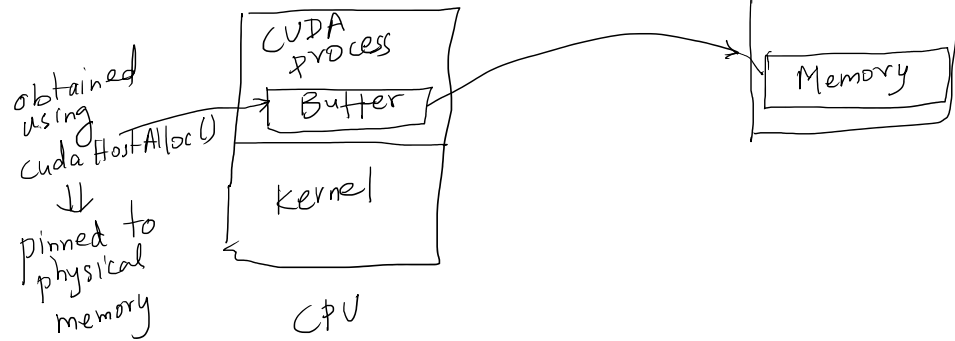· dest    address (device)
· # of  bytes to transfer

$\rightarrow$ DMA happens in background
$\rightarrow$ Frees up CPU for other jobs.

· DMA transfer is a two step process

CPU

user space { CUDA process
[ 4 ]   ①

kernel space { Driver buffer

GPU

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
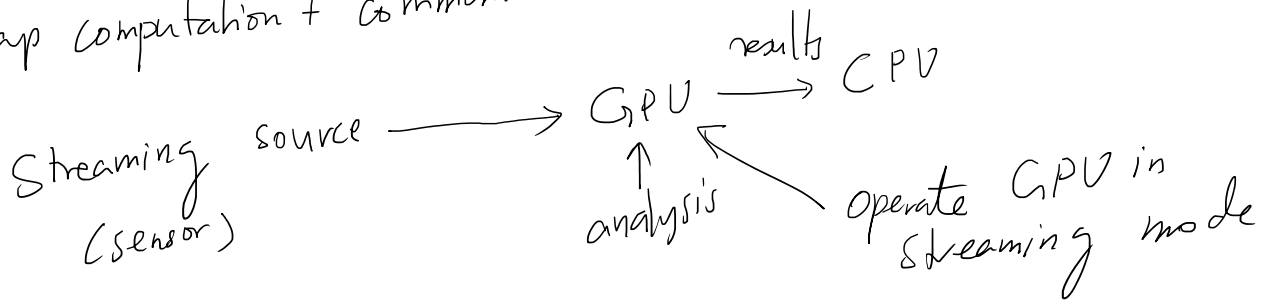
Code example: *streams*

So far: data is transferred in full
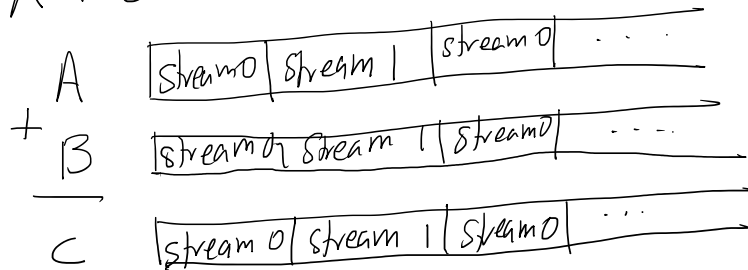before computation begins

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
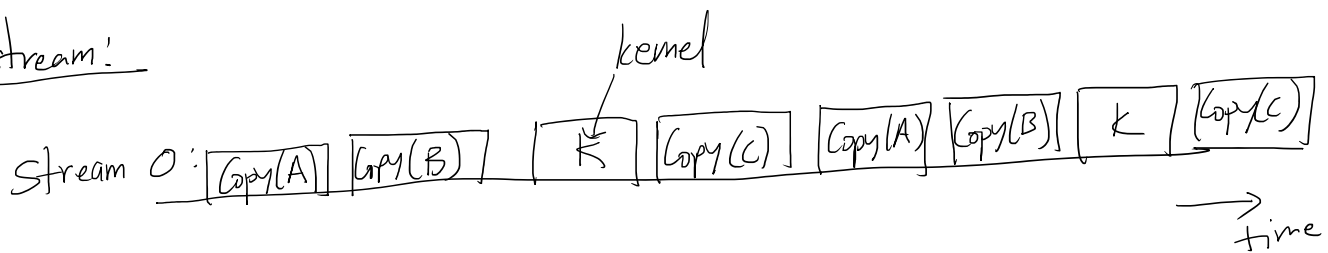4) Overlap computation + communication

Streaming source $\longrightarrow$ GPU $\xrightarrow{\text{results}}$ CPU
(sensor)

analysis $\uparrow$ $\nwarrow$ Operate GPU in
streaming mode

Vector addition:

$$C = A + B$$

$+$ A | Stream 0 | Stream 1 | Stream 0 | . . .

$+$ B | Stream 0 | Stream 1 | Stream 0 | . . .

‾‾‾

C | 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:

Compute Engine — executes kernel

Transfer Engine

Transfer Engine

TE0 : (host → device)

TE1 : (device → host)

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

Stream 1 :   Copy(A)   Copy(B)   K   Copy(C)   Copy(A)   ....

TE0    TE0    TE1    TE0    TE0

TE0    TE0    TE1

time →

overlapped computation and communication
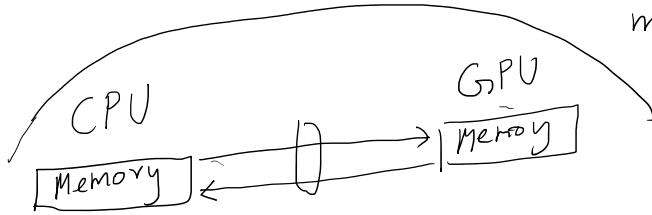
· All operations are done in async fashion using pinned memory.

# Unified memory

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

Page table
valid dirty

Launch kernel

GPU

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

CPU — A: { page 0, page 1, page 2, page 3 }

Page table:
| | valid | dirty |
|------|-------|-------|
| page 0 | 1 | 0 |
| page 1 | 1 | 0 |
| page 2 | ~~1~~ 0 | ~~1~~ 0 |
| page 3 | 1 | 0 |

GPU:
| | V | D |
|------|---|---|
| page 0 | 0 | 0 |
| page 1 | 0 | 0 |
| page 2 | ~~0~~ 1 | 0 |
| page 3 | 0 | 0 |

Code example: *cuda_omp*

OpenMP threads on CPU

$C = A + B$

$S_0$    $S_1$    $S_2$

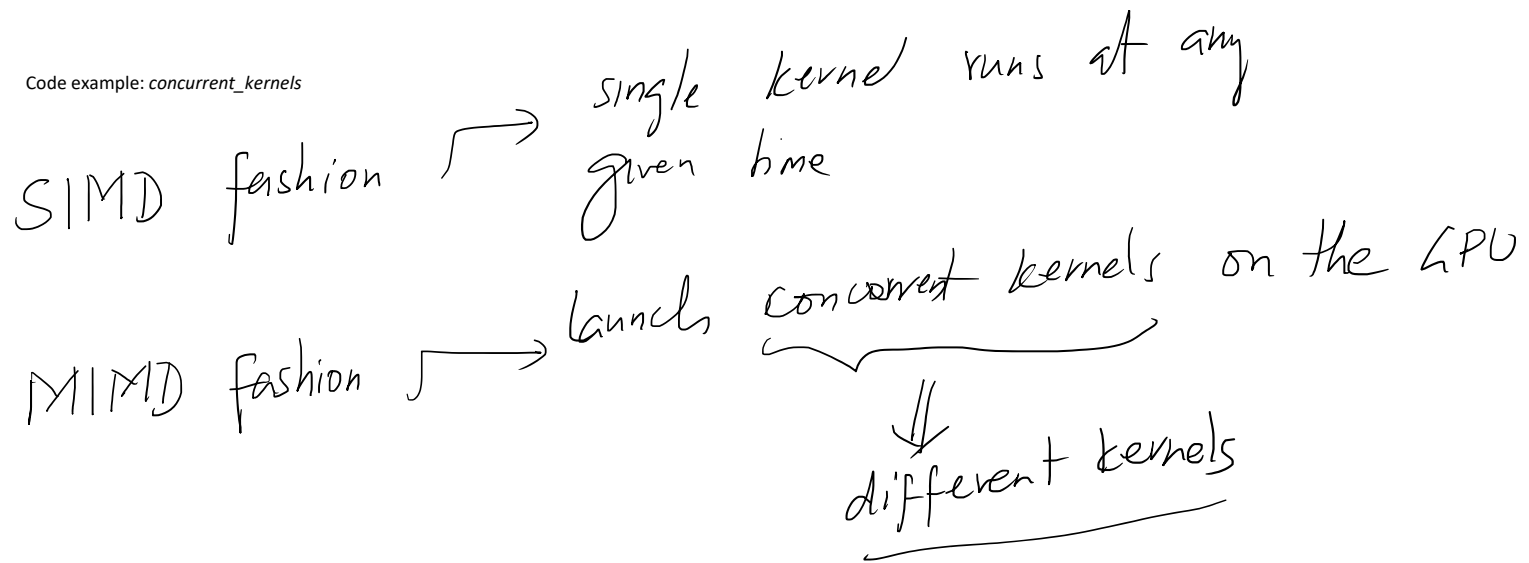- Extract coarse-grained parallelism using OpenMP

A

+

B

||

C

$SSSS...S$

↑ CUDA threads per openMP thread

. Extract fine-grained parallelism using CUDA

# Concurrent kernel execution

Code example: *concurrent_kernels*

SIMD fashion $\longrightarrow$ single kernel runs at any given time

MIMD fashion $\longrightarrow$ 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

- *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