

Exercise 6 – GPUs

S. Kuckuk, D. Ernst, G. Wellein, G. Hager

Erlangen National High Performance Computing Center (NHR@FAU)

Assignment 6 – Task 1

Performance: $P_{chip} = n_{core} \cdot n_{super}^{FP} \cdot n_{FMA} \cdot n_{SIMD} \cdot f$

#Cores Super-scalarity FMA factor SIMD factor Clock Speed

- For one Alex node, FP64

- $P_{CPU} = 2 \cdot 64 \text{ cores} \cdot 2 \frac{\text{instruction}}{\text{cycle}} \cdot 2 \frac{\text{Flop}}{\text{FMA}} \cdot 4 \frac{\text{FMA}}{\text{instruction}} \cdot 2 \text{ Ghz} = 4096 \frac{\text{GFlop}}{\text{s}}$
- $P_{GPU} = 108 \text{ SM} \cdot 1 \frac{\text{instruction}}{\text{cycle}} \cdot 2 \frac{\text{Flop}}{\text{FMA}} \cdot 32 \frac{\text{FMA}}{\text{SM}} \cdot 1.41 \text{ Ghz} = 9746 \frac{\text{GFlop}}{\text{s}}$
- $P_{8-GPU} = 8 \cdot P_{GPU} = 78 \frac{\text{TLop}}{\text{s}}$

Assignment 6 – Task 2

- #FLOP for a **D**ouble **G**eneral **P**recision **M**atrix **M**atrix Multiplication (DGEMM) for matrices with dimensions $MK \times KN = MN$

- $N_{flops} = 2 \cdot MNK$

- Assuming compute boundness

- $T_{CPU} = \frac{N_{Flops}}{P_{CPU}} = \frac{2MNK}{P_{CPU}} = \frac{2 \cdot 2048^3}{4096 \text{ GFlop/s}} = 4194 \mu s$

- $T_{8-GPU} = \frac{N_{Flops}}{P_{8-GPU}} = \frac{2MNK}{P_{8-GPU}} = \frac{2 \cdot 2048^3}{78 \text{ TFlop/s}} = 220 \mu s$

Assignment 6 – Task 2

- Assuming compute boundness

- $$T_{CPU} = \frac{N_{Flops}}{P_{CPU}} = \frac{2MNK}{P_{CPU}} = \frac{2 \cdot 2048^3}{4096 \text{ GFlop/s}} = 4194 \mu\text{s}$$

- $$T_{8-GPU} = \frac{N_{Flops}}{P_{8-GPU}} = \frac{2MNK}{P_{8-GPU}} = \frac{2 \cdot 2048^3}{78 \text{ TFlop/s}} = 220 \mu\text{s}$$

- Transfer times

- Two shared PCIe links per node, each with 25GB/s per direction
- Host-to-device: 2 input matrices, device-to-host: one result matrix

- $$T_{PCIe} = \frac{3 \cdot V_{matrix}}{2 \cdot B_{PCIe}} = \frac{3 \cdot M^2 \cdot 8 \text{ B}}{2 \cdot B_{PCIe}} = \frac{3 \cdot 2048^2 \cdot 8 \text{ B}}{2 \cdot 25 \frac{\text{GB}}{\text{s}}} = 2013 \mu\text{s}$$

Assignment 6 – Task 2

- Putting it together

- $T_{CPU} = \frac{N_{Flops}}{P_{CPU}} = \frac{2MNK}{P_{CPU}} = \frac{2 \cdot 2048^3}{4096 \text{ GFlop/s}} = 4194 \mu s$

- $T_{8-GPU} = \frac{N_{Flops}}{P_{8-GPU}} = \frac{2MNK}{P_{8-GPU}} = \frac{2 \cdot 2048^3}{78 \text{ TFlop/s}} = 220 \mu s$

- $T_{PCIe} = \frac{3 \cdot V_{matrix}}{2 \cdot B_{PCIe}} = \frac{3 \cdot M^2 \cdot 8 B}{2 \cdot B_{PCIe}} = \frac{3 \cdot 2048^2 \cdot 8 B}{2 \cdot 25 \frac{GB}{s}} = 2013 \mu s$

- $S = \frac{T_{CPU}}{T_{8-GPU} + T_{PCIe}} = \frac{4194 \mu s}{220 \mu s + 2013 \mu s} = 1.88$

Assignment 6 – Task 3

- For one Frontier node, FP64

- $P_{CPU} = 64 \text{ cores} \cdot 2 \frac{\text{instruction}}{\text{cycle}} \cdot 2 \frac{\text{Flop}}{\text{FMA}} \cdot 4 \frac{\text{FMA}}{\text{instruction}} \cdot 2 \text{ Ghz} = 2048 \frac{\text{GFlop}}{\text{s}}$

- $P_{GPU} = 110 \text{ CU} \cdot 1 \frac{\text{instruction}}{\text{cycle}} \cdot 2 \frac{\text{Flop}}{\text{FMA}} \cdot 64 \frac{\text{FMA}}{\text{CU}} \cdot 1.7 \text{ Ghz} = 24 \frac{\text{TFlop}}{\text{s}}$

- $P_{8-GPU} = 8 \cdot P_{GPU} = 192 \frac{\text{TFlop}}{\text{s}}$

Assignment 6 – Task 4

■ DGEMM on Frontier

- $T_{CPU} = \frac{N_{Flops}}{P_{CPU}} = \frac{2MNK}{P_{CPU}} = \frac{2 \cdot 2048^3}{2048 \text{ GFlop/s}} = 8389 \mu\text{s}$
- $T_{8-GPU} = \frac{N_{Flops}}{P_{8-GPU}} = \frac{2MNK}{P_{8-GPU}} = \frac{2 \cdot 2048^3}{192 \text{ TFlop/s}} = 89.5 \mu\text{s}$
- $T_{Link} = \frac{3 \cdot V_{matrix}}{8 \cdot B_{Link}} = \frac{3 \cdot M^2 \cdot 8 \text{ B}}{8 \cdot B_{Link}} = \frac{3 \cdot 2048^2 \cdot 8 \text{ B}}{8 \cdot 36 \frac{\text{GB}}{\text{s}}} = 350 \mu\text{s}$

■ Speed-up CPU vs GPU

- $S = \frac{T_{CPU}}{T_{8-GPU} + T_{Link}} = \frac{8389 \mu\text{s}}{89.5 \mu\text{s} + 350 \mu\text{s}} = 19$

Assignment 6 – Task 4

- Speed-up Alex vs Frontier

- $S = \frac{T_{Alex}}{T_{Frontier}} = \frac{2233 \mu s}{440 \mu s} = 5$

Assignment 6 – Task 5

- Little's Law

$$L = \frac{N}{b_s} \leftrightarrow N = L \cdot b_s$$

The diagram illustrates the relationship between the variables in Little's Law. The equation $L = \frac{N}{b_s} \leftrightarrow N = L \cdot b_s$ is shown in blue. Three callout boxes are connected to the variables: a box labeled "#Bytes in flight" points to N , a box labeled "Latency" points to L , and a box labeled "Bandwidth" points to b_s .

Assignment 6 – Task 5

- Kernel code (slightly modified)

```
__global__ void kernel(double* A, double* B, double* C, int size) {  
    int x = threadIdx.x + blockDim.x * blockIdx.x;  
    int y = threadIdx.y + blockDim.y * blockIdx.y;  
  
    if (x >= size || y >= size) return;  
  
    C[x + y * size] += A[x + y * size] * B[y];  
}
```

Two DRAM loads and
one DRAM store
(assuming B fits into L2)

Assignment 6 – Task 5

- Little's Law

- $L = \frac{N}{b_S} \leftrightarrow N = L \cdot b_S$

- $L_{seconds} = \frac{L_{cycles}}{F}$

- $N_{bytes} = N_{threads} \cdot V_{thread}$

- $V_{thread} = 3 \cdot 8 B = 24 B$

- Put together

- $N_{threads} = \frac{670 \text{ cycles}}{1.41 \text{ GHz}} \cdot \frac{2039 \frac{GB}{s}}{24 \frac{B}{thread}} = 40370 \text{ threads}$

Minimum number of threads required to saturate the DRAM bandwidth, actual number is higher

Assignment 6 – Task 5

- Kernel code (slightly modified)

```
__global__ void kernel(double* A, double* B, double* C, int size) {  
    int x = threadIdx.x + blockDim.x * blockIdx.x;  
    int y = threadIdx.y + blockDim.y * blockIdx.y;  
  
    if (x >= size || y >= size) return;  
  
    C[x + y * size] += A[x + y * size] * B[y];  
}
```

Threads with same y can
share data in L1
-> blocks with large x extent
reduce L2 cache traffic

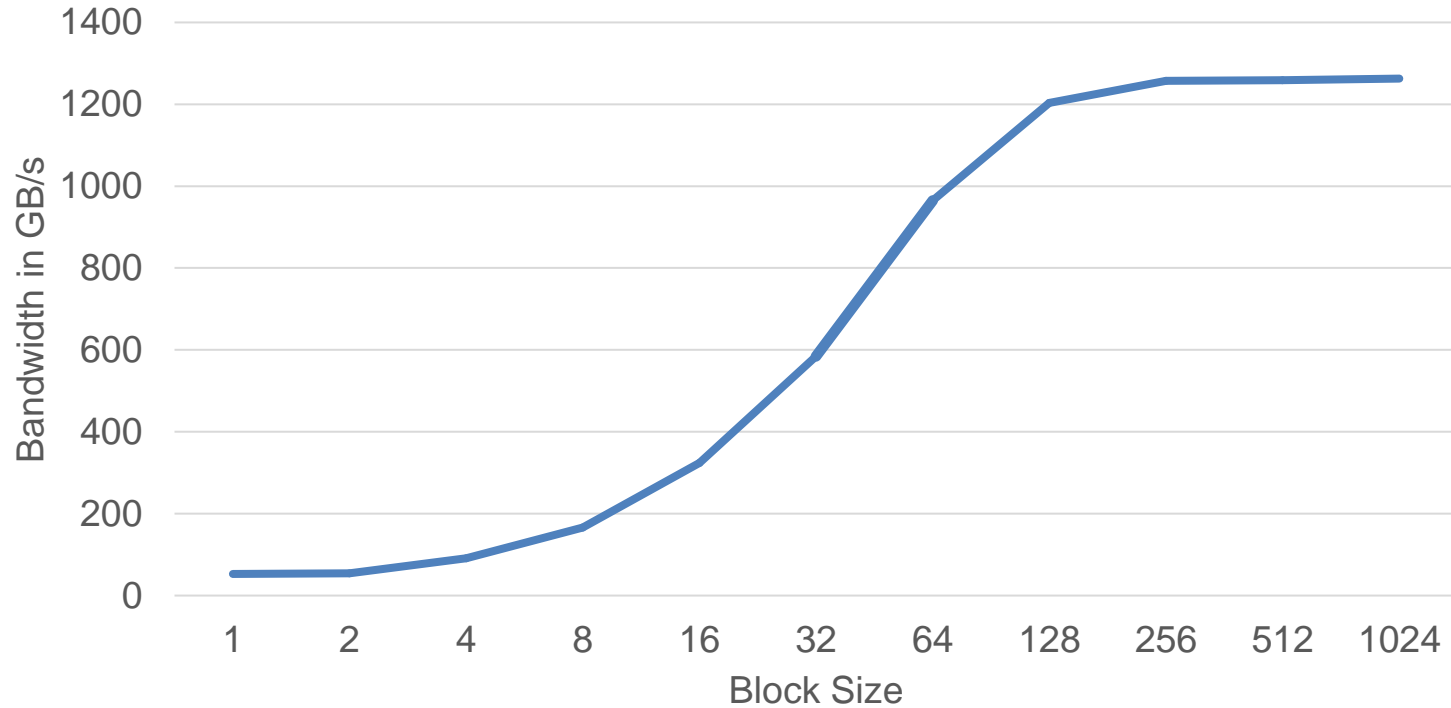
Assignment 6 – Task 6

- Kernel code for $a[i] = a[i] + 1.0$

```
__global__ void kernel(double* a, int size) {  
    int x = threadIdx.x + blockDim.x * blockIdx.x;  
    int stride = blockDim.x * blockDim.x;  
  
    for ( ; x <= size; x += stride)  
        a[x] += 1;  
}
```

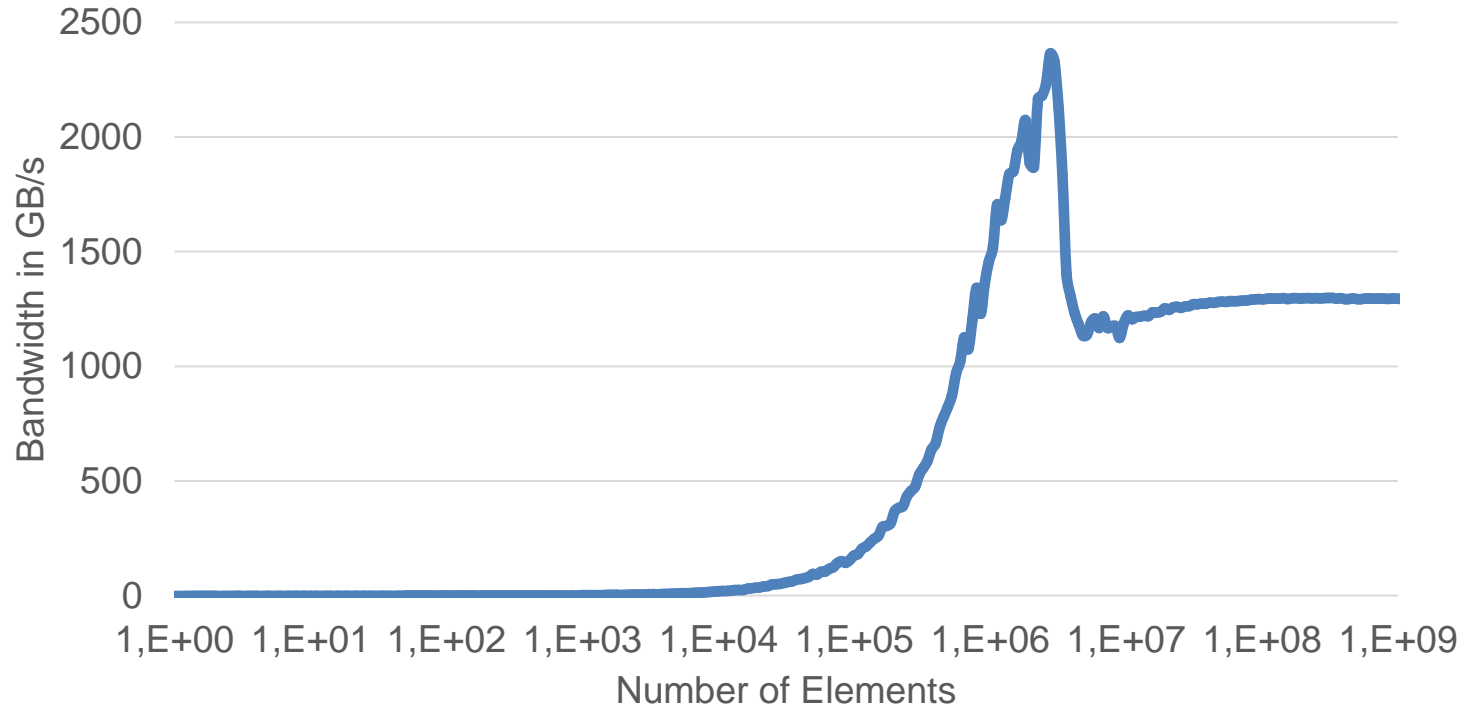
Assignment 6 – Task 6

- Block size investigation



Assignment 6 – Task 6

- DRAM Bandwidth



Assignment 6 – Task 6

■ PCIe Bandwidth

