

作业(2) : DLP & GPU

截至时间 : 2022.11.25/周五 23:59:59

提交方式 : 超算习堂 (<https://easyhpc.net/course/157>)

Q1-Sol:

Peak throughput = $1.5 \times 16 \times 16 = 384$ GFLOPS/s (or, $1.5 \times 16 \times 16 \times 2 = 768$ GFLOPS/s)

Assume each single-precision operation requires 2 four-byte operands and outputs 1 four-byte result, sustaining this throughput would require $12 \text{ bytes/FLOP} \times 384 \text{ GFLOPS/s} = 4.6 \text{ TB/s}$ of bandwidth. As such, the throughput is not sustainable on the 100 GB/s off-chip bandwidth, but can still be achieved in short bursts when using on-chip memory.

Q2-Sol:

- $4 \times (2048 / 16) = 4 \times 128 = 512$
- $16\text{KB} / 4\text{KB} = 4$
- register spilling, leading to use local memory which is likely to downgrade performance.
- 50%, only 16 blocks can be supported. $16 \times 16 / 512 = 50\%$
- $\langle\langle\langle 16, 32 \rangle\rangle\rangle, \langle\langle\langle 8, 64 \rangle\rangle\rangle, \langle\langle\langle 4, 128 \rangle\rangle\rangle$

Q3-Sol:

Any differences on core/instruction/link are acceptable. Here are some examples:

- new 8-bit FP8 floating point format on H100 tensor core
- new DPX instructions to accelerate dynamic programming
- distributed shared memory
- larger L2 cache
- new thread block cluster
- HBM3 memory system
- Faster NVLink
-

Q4-Sol:

N.A.