Question 1:

How long does it take to load 200B bf16 model into the systolic array of 32 TPU v4p from HBM

TPU v4p HBM to Tensor Core bandwidth

$$= 1.2 \times 10^{12}$$

Model size $= 2 \times 10^9$

$$\text{Sharded\_model} = \frac{2 \times 2 \times 10^9 \times 100}{32}$$

$$\frac{\text{Sharded\_model}}{\text{HBM to Tensorcore BW}} = \frac{\frac{2 \times 2 \times 10^9 \times 100}{32 \cdot 16}}{1.2 \times 10^{12}} = \frac{2 \times 1 \times 100}{1.2 \times 16 \times 10^3} \sec$$

$$= \frac{2 \times 1 \times 100}{19.2 \times 10^3}$$

$$= 0.052 \times 10^{-3} \times 100$$

$$= 100 \times 2 \times 5.2 \times 10^{-5} \sec$$

$$= 10.4 \times 10^{-5} \times 100$$
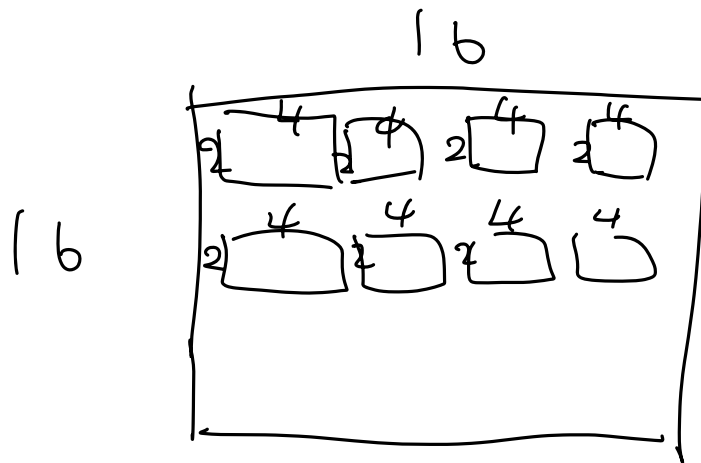
$$= 10.4 \times 10^{-3}$$

$$= 10 \text{ ms}$$

Question 2: TPU v5e pod

① How many CPU hosts are there?

TPU v5e Pod size = 16 × 16

Every 4×2 contains a CPU

So = $\frac{16}{4} \times \frac{16}{2}$ = 4 × 8 = 32 CPU hosts

16

16

② How many TPU Tensor Cores?

16 × 16 × 1 = 256

TPU v5e has TPU core per chip

③ Total FLOPs/sec for the whole pod

$1.97 \times 10^{14} \times 16 \times 16 = 504.32 \times 10^{14} = 5.04 \times 10^{16}$

④ Total HBM = $16\,GB \times 16 \times 16$

$$= 4096\,GB = 4\,TB$$

Question 3: PCIe operational intensity

$$\boxed{\begin{array}{l} A = bfloat[D, F] \\ X\,Activations = bfloat[B, D] \end{array}}$$ on DRAMs

$$FLOP_s = BF(\underset{mul}{\underline{D}} + \underset{add}{\underline{D-1}})$$

$$= BF(2D-1) = 2BD\bar{F} - BF$$

$$\simeq 2BD\bar{F} \ \checkmark$$

$$FLOP_s = 2BDF \overset{F=4D}{=} 2BD \cdot 4D = 8BD^2$$

$$Bytes = 2\underset{Load}{(\underline{DF + BD})} + \underset{write}{\underline{2BF}}$$

$$\overset{F=4D}{=} 2(4D^2 + BD) = 8D^2 + 2BD \simeq 8D^2$$

$$\frac{FLOP_s}{Bytes} = \frac{8BD^2}{8D^2} = B$$

$$AI_{hardware} = \frac{\cancel{4\times 2} \times 9.2 \times 10^{14}}{1.5 \times 10^{10}}$$

Assume single chip

$$= \frac{8 \times 9.2 \times 10^4}{1.5} = \cancel{4.9 \times 10^4} \quad 61250$$

$$= 4.9 \times 10^5$$

Question 4: general matmul latency

Multiply $int8[16384, 4096]$ by $int8[B, 4096]$

Assumes: 1 TPU v5e (chip)

1. ✓

Time $= \max(\text{bytes load/write}, \text{FLOPs time})$

$$= \max\left( \frac{16384 \times 4096 + B \cdot 4096 + B \cdot 16384}{8.1 \times 10^{11}}, \right.$$

$$\left. \frac{B \times 16384 \times 2 \times 4096}{3.94 \times 10^{14}} \right)$$

$$= \max\left( \frac{6.7 \times 10^7 + 20480 B}{8.1 \times 10^{11}} + \frac{1.3 \times 10^8 B}{3.94 \times 10^{14}} \right)$$

Compute-bound when
$$\frac{1.3 \times 10^8 B}{3.94 \times 10^{14}} > \frac{6.7 \times 10^7 + 2 \times 10^4 B}{8.1 \times 10^{11}}$$

$$\Longrightarrow \quad \frac{1.3 B}{3.94 \times 10^6} > \frac{6.7 \times 10^3 + 2B}{8.1 \times 10^7}$$

$$\Longleftrightarrow \quad 1.3 \times 8.1 \times 10^7 B > 6.7 \times 10^3 \times 3.94 \times 10^6$$
$$+ 2 \times 3.94 \times 10^6 B$$

$$\Longleftrightarrow \quad 10.53 \times 10^7 B - 7.88 \times 10^6 B > 26.398 \times 10^9$$

$$\Longleftrightarrow \quad 105.3 \times 10^6 B - 7.88 \times 10^6 B > 2.6 \times 10^{10}$$

$$\Longleftrightarrow \quad 97.42 \times 10^6 B > 2.6 \times 10^{10}$$

$$\Longleftrightarrow \quad B > \frac{2.6 \times 10^4}{97.42}$$

$$\Longleftrightarrow \quad B > 268$$

2. Assumes: bandwith $= 22 \times 8.1 \times 10^{11}$

Compute-bound when:
$$\frac{1.3 \times 10^8 B}{3.94 \times 10^{14}} > \frac{6.7 \times 10^7 + 2 \times 10^4 B}{22 \times 8.1 \times 10^{11}}$$

$$\Longrightarrow \quad \frac{1.3 \times B}{3.94 \times 10^6} \quad {}^{) 10^8} \quad \underset{}{\overset{) 10^{-4}}{>}} \quad \frac{6.7 \times 10^3 + 2 B}{22 \times 8.1 \times 10^7}$$

$\iff$  $1.3 \times 22 \times 8.1 \times 10^7 B > 3.94 \times 6.7 \times 10^9$

$\qquad\qquad\qquad + 2 \times 3.94 \times 10^6 B$

$\iff$  $231.66 \times 10^7 B > 26.398 \times 10^9$

$\qquad\qquad\qquad + 7.88 \times 10^6 B$

$\iff$  $2316.6 \times 10^6 B - 7.88 \times 10^6 B > 2.6398 \times 10^{10}$

$\iff$  $2308.72 \times 10^6 B > 2.6398 \times 10^{10}$

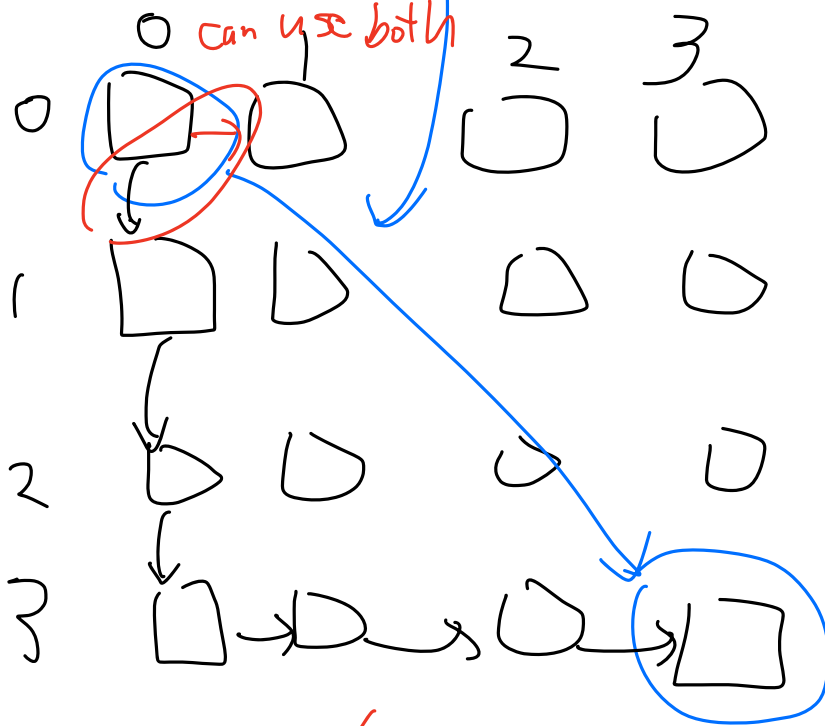$\iff$  $2.30872 \times 10^9 B > 2.6398 \times 10^{10}$

$\iff$  $B > \dfrac{2.6398 \times 10^{10}}{2.30872 \times 10^9}$

$\qquad\qquad = \underline{11.43}$  \#

---

# Question 5: ICI bandwidth

TPU v5e : 4x4 slice

bfloat [8, 128, 8192]

can use both

0   1   2   3

1. $6 \times 1\mu s$ = first byte latency
   (hops)

2. Total time $= 6\mu s + \dfrac{2 \times 8 \times 128 \times 8192}{2 \times 4.5 \times 10^{10}}$

   $= 6\mu s + \dfrac{16,777,216}{2 \times 4.5 \times 10^{10}}$

   $= 6\mu s + \dfrac{1.6 \times 10^7}{2 \times 4.5 \times 10^{10}}$

   $= 6\mu s + \dfrac{1.6}{2 \times 4.5 \times 10^3}$

   $= 6 \times 10^{-6} sec + 0.35 \times 10^{-3} sec$

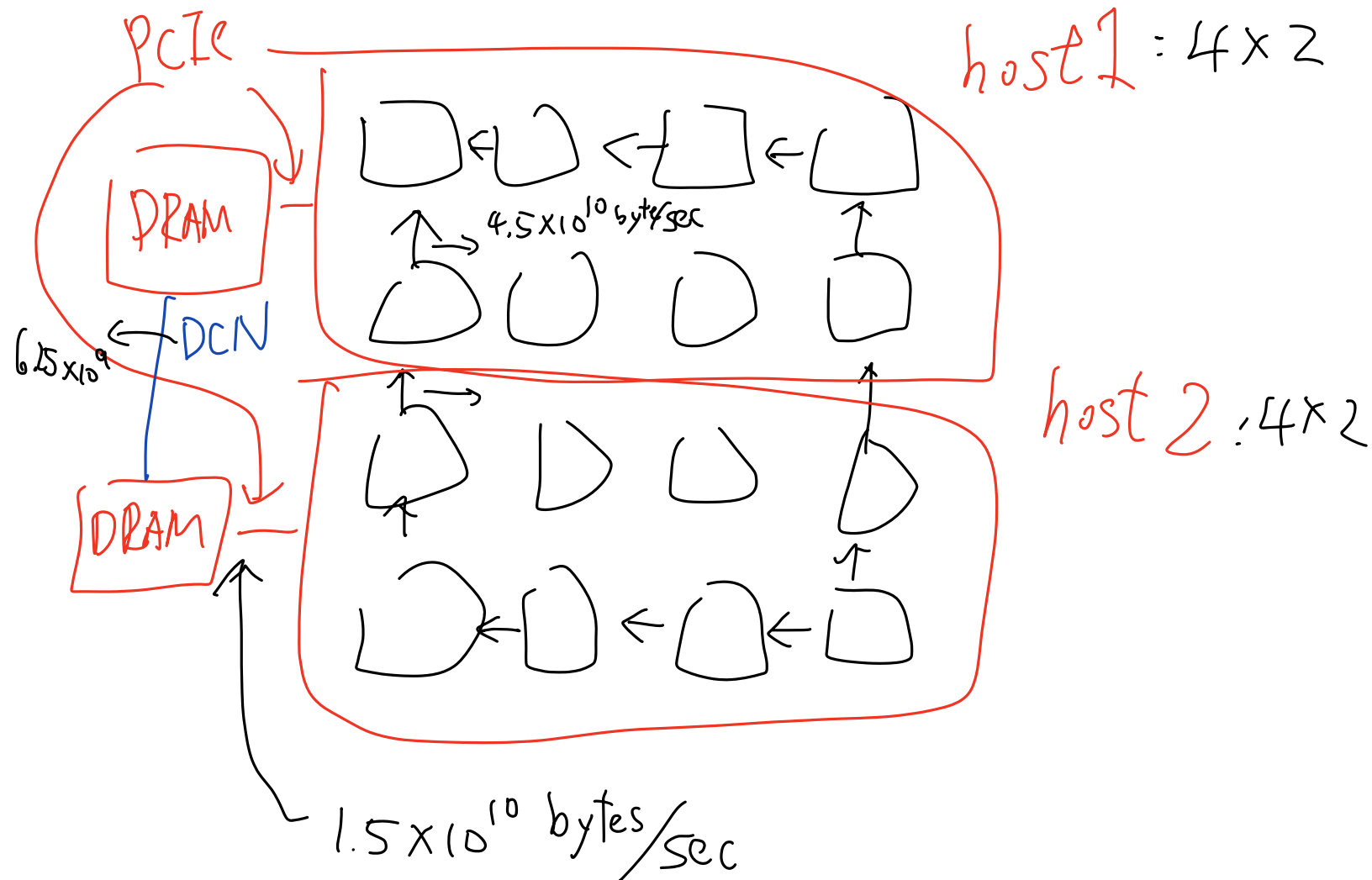   $= 6 \times 10^{-6} sec + 3.5 \times 10^{-4} sec$

$$= 3\text{-}5\text{-}6 \ \mu s$$
181

Question 6:

Matrix A: int8 [128×1024, 128×1024]

Sharded evenly across TPU v5e 4×4 slice
but they are on host's DRAM on each
chip.



PCIe

DRAM

DCN

6.5×10⁹  →  $6.5 \times 10^9$

DRAM

host 1: 4×2

4.5×10¹⁰ bytes/sec  →  $4.5 \times 10^{10}$ bytes/sec

host 2: 4×2

1.5×10¹⁰ bytes/sec  →  $1.5 \times 10^{10}$ bytes/sec

First: We can choose either:

① Copy half of the data through DCN to host 1's DRAM and then copy things to HBM from DRAM

or ② Each host copy its data to TPU chips and then utilize ICI to copy data

Option ① : DCN throughput is slow, so use ②

Steps : (will take the max over those numbers)

1. Copy data from DRAM to HBM

2. Transfer data from all chips to TPU {0, 0}        $2^7$   $2^{10}$

3. Compute int8 [ $\overset{||}{128} \times \overset{||}{1024}$, $128 \times 1024$]
   bfloat16 [ $\overset{||}{8}$ , $128 \times 1024$]

**1.**

$$128 \times 1024 \times 128 \times 1024 = 2^{34} = 2^4 \text{ GB}$$
$$= 16 \text{ GiB}$$

$$\frac{16 \text{ GiB}}{2} = 8 \text{ GiB} = \text{data on each host}$$

Time to load from DRAM to $\underline{\text{HBM}}$ <span style="color:red">to all TPU chips through PCIe</span>

$$= \frac{8 \text{ GiB} \checkmark}{16 \times 1.5 \times 10^{10} \text{ bytes/sec}} \quad \color{red}= 0.035 \text{ sec} = 35 \text{ ms}$$

**2.** Time to send data from HBM to TPU$\{0,0\}$ through ICI :

<span style="color:red">15 GiB because TPU$\{0,0\}$ needs to receive</span>

$$\frac{15\cancel{8} \text{ GiB}}{2 \times 4.5 \times 10^{10}} = \frac{\overset{15}{\cancel{8}} \text{GiB} \quad \color{red}{15 \text{ GiB}}}{9 \times 10^{10} \text{ bytes/sec}}$$

$$\color{red}= 0.178 \text{ sec} = 178 \text{ ms}$$

**3.** Time to load and write <span style="color:red">to MXU</span>

Load int8 $[128 \times 1024, 128 \times 1024] = \color{red}\dfrac{128 \times 1024 \times 128 \times 1024}{8.1 \times 10^{11}}$
<span style="color:red">$= 0.021$</span>
<span style="color:red">$= 21 \text{ ms}$</span>

Load bf16 $[8, 128 \times 1024] =$ I think the problem assume it is already in TPU HBM

Write bf16 $[8, 128 \times 1024] = 2 \times (8 \times 128 \times 1024)$
$$= 2 \times 2^3 \times 2^7 \times 10^{10}$$

$$\frac{2\text{MiB}}{8.1 \times 10^{11} \text{ (HBM speed)}} = \text{neglibie} \qquad = 2\text{MiB}$$

3. Compute time :

$$\frac{2BDF}{1.97 \times 10^{14} \text{ FLOP/sec}} = \frac{2 \times (8 \times 128 \times 1024 \times 128 \times 1024)}{1.97 \times 10^{14} \text{ FLOPs/sec}}$$

$$= \frac{2 \times (2^3 \times 2^7 \times 2^{10} \times 2^7 \times 2^{10})}{1.97 \times 10^{14}}$$

$$= \frac{2^{38}}{1.97 \times 10^{14}}$$

$$= \frac{256 \text{ GiB}}{1.97 \times 10^{14}} \qquad = 0.00139 \text{ se}$$

$$= 1.3 \text{ ms}$$

Total = max stage

$$= 178 \text{ ms}$$