$$T_{comm} = \frac{\text{Comm bytes}}{\text{bandwith bytes/sec}}$$

depend on how you write code

Constant, determined by hardware

$$T_{math} = \frac{\text{Computation FLOPs}}{\text{Accelerator FLOPs/sec}}$$

$$T_{lower} = \max(T_{math}, T_{comms})$$

( assuming overlap computation)

$$\text{Arithmetic intensity} = \frac{\text{FLOPs}}{\text{bytes}}$$

When

$T_{math} > T_{comm}$, this means communication can be overlapped with $T_{math}$,

so actual $\frac{\text{FLOPs}}{\text{sec}} = \frac{\text{Computation FLOPs}}{T_{math}}$

$$= \text{Accelerator} \frac{\text{FLOPs}}{\text{sec}}$$
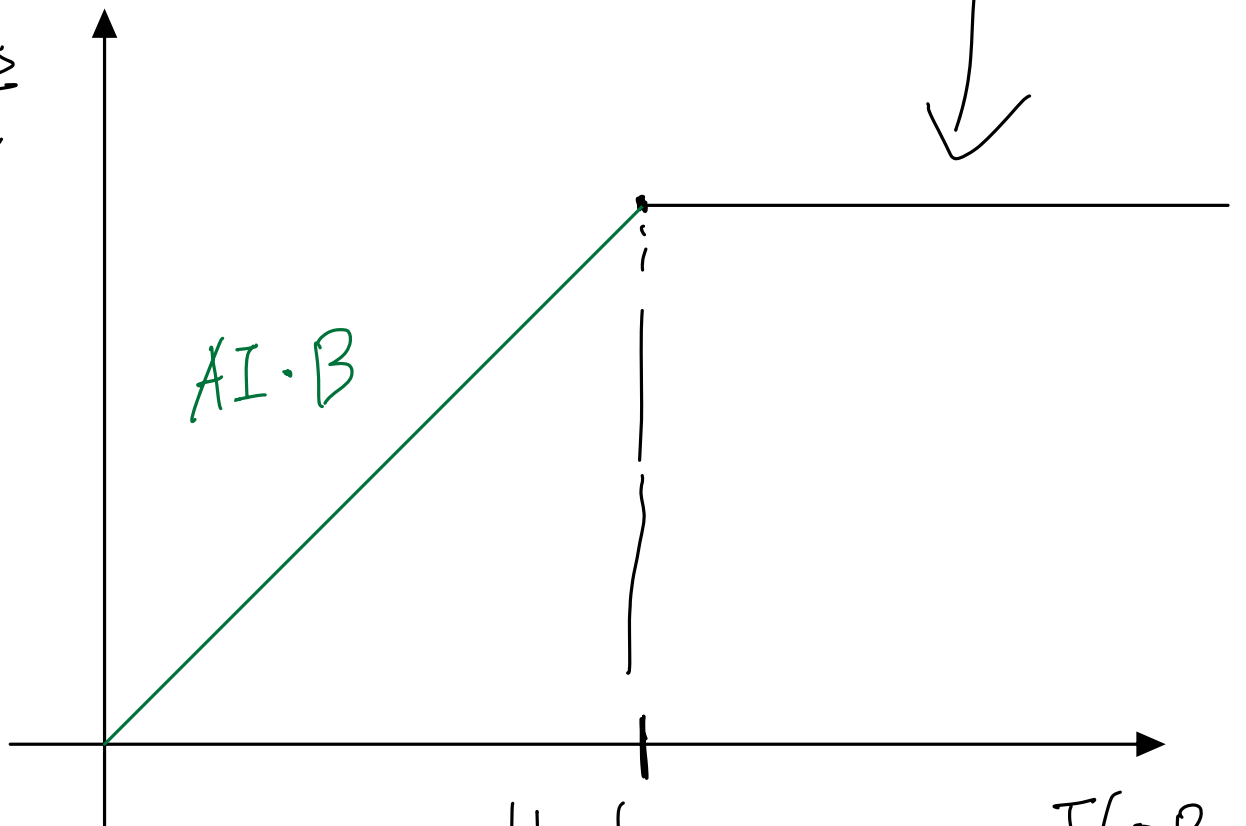
( constant)

$T_{math} > T_{comm}$ also means

$(\equiv)$ $\dfrac{\text{computation FLOPs}}{\text{Accelerator FLOPs}} > \dfrac{\text{Communication bytes}}{\text{bandwith bytes}/sec}$

$(\equiv)$ Arithmetic Intensity $>$ AI $_{(hardware)}$
( compute)

Summary: when AI compute $>$ AI hardware

actual $\dfrac{\text{FLOPs}}{sec}$ realized is

always Accelerator FLOS$/sec$

When

$$T_{math} < T_{comm}$$

$$\Longleftrightarrow AI_{compute} < AI_{hardware}$$

This also means the actual throughput is

$$\frac{FLOPs}{T_{comm}} = \frac{FLOPs}{\left(\frac{\text{Communication bytes}}{\text{bandwith bytes}/\text{sec}}\right)}$$

$$= \frac{FLOPs}{\text{Communication bytes}} \cdot \left(\text{bandwith bytes}/\text{sec}\right)$$

$$= AI_{compute} \cdot B$$

This means, as $AI_{compute}$ increase

the realized $\frac{FLOPs}{sec}$ increases linearly

Combining results from:

① $T_{math} < T_{comm}$

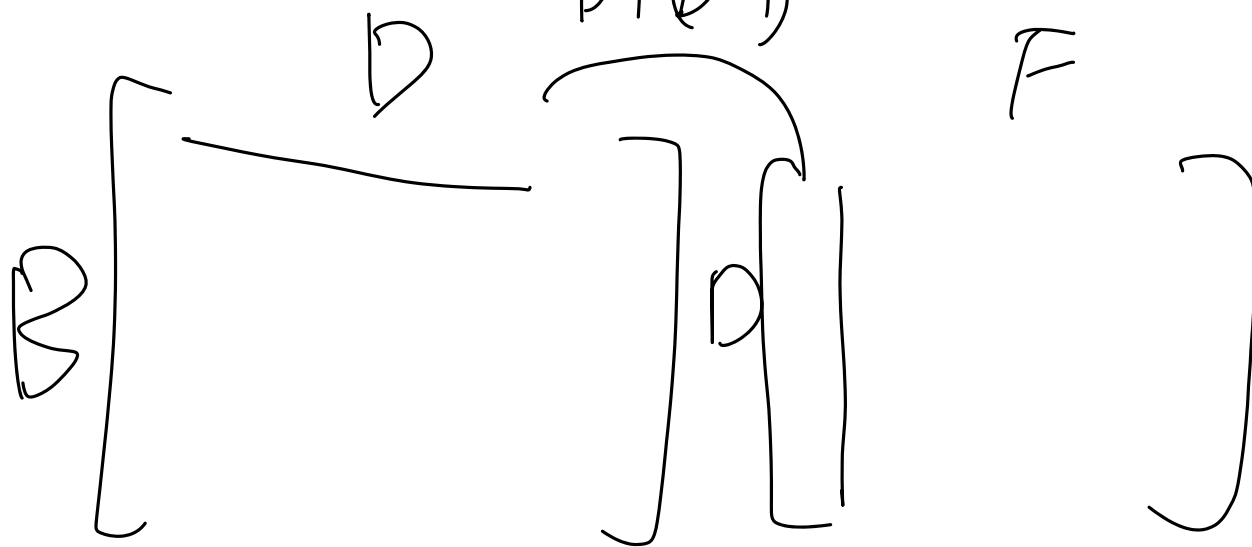② $T_{math} > T_{comm}$

observations, we get the roofline model

# Matrix multiplication

$$X = (B, D) \text{ bf16} \rightarrow 16\text{bit} \rightarrow 2\text{byte}$$
$$Y = (D, F) \text{ bf16} \xrightarrow{\quad} (B, \bar{F})$$

$$* $$
$$D+(D-1) \leftarrow +$$



$$\text{Load} = 2\text{byte} \cdot B \cdot D + 2 \cdot D \cdot F$$

$$\#\text{FLOPS} = (2D-1) \cdot B \cdot F$$
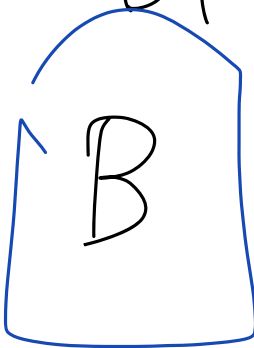
$$\sim = 2BDF$$

$$\text{Write} = 2\underset{\text{bytes}}{B} F$$

$$\text{Intensity} = \frac{2BDF}{2BD + 2DF + 2BF}$$
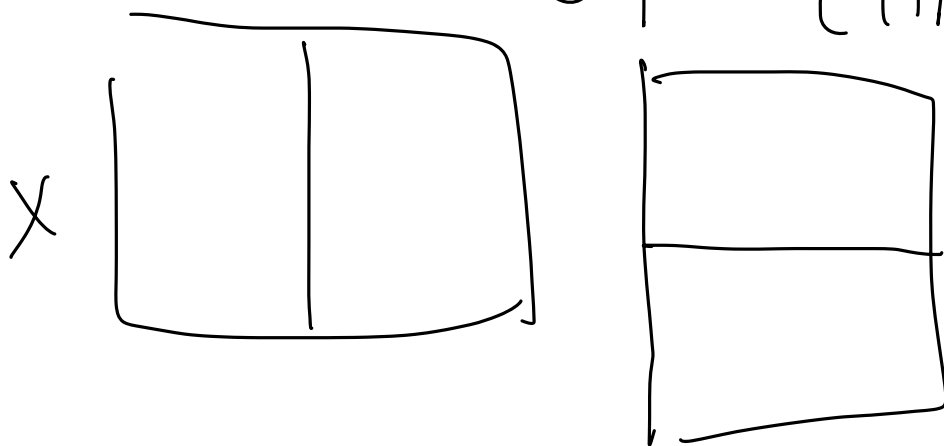
$$= \frac{BDF}{BD + DF + BF}$$

$$\approx \frac{BDF}{DF}$$

$$= \boxed{B}$$

(assume D, F are dominating in Transformer cases which is hidden dimensions)

# Network communication rooflines

$X: \text{bfloat16}[B, D]$    matmul
$Y: \text{bfloat16}[D, F]$    $\rightarrow \text{bfloat}[B, F]$

across 2 TPUs

Consider:   $X[:, :D//2] @ Y[:D//2, :]$
$+ \; X[:, D//2:] @ Y[D//2:, :]$
$= \; X @ Y$    (tiling)



X

Assume:   Network bandwidth $= 4.5 \times 10^{10}$ bytes/sec
    FLOPs/sec each chip $= 1.97 \times 10^{14}$ FLOPs/sec

$$T_{math} = \frac{2 BDF}{2 \cdot 1.9 \times 10^{14}} \; \div \; \frac{BDF}{1.9 \times 10^{14}}$$

$\uparrow$
2 chips

Assume the best case when we fully utilize a

$$T_{comms} = \frac{BF + BF}{4.5 \times 10^{10}}$$ (left TPU send BF to right <span style="color:red">Single TPU</span> and right send BF to left)

When

$$T_{math} > T_{comms} \text{ happens}$$

$$(\Longleftrightarrow) \quad \frac{D}{2} = \frac{BDF}{2BF} > \underbrace{\frac{1.97 \times 10^{14}}{4.5 \times 10^{10}}}_{\text{Arithmetic intensity}}$$

$$(\Longrightarrow) \quad D > 8755$$

$(\Longleftrightarrow)$ When $D > 8755$, we know that we know 2 TPU is compute-bound (means good), assuming each TPU is fully utilizing its FLOPs/sec (the best case)

$\Rightarrow$ so this means the previous B needs to be $> 240$ (hardware

# Question 1:

$$X[B,D] \,@\, Y[D,F] \to Z[B,F] \text{ in}$$

int 8 ( 1 byte) instead of bfloat16

1. Loaded = $BD + DF$

   Write = $BF$

2. $\text{FLOPs} = B \cdot F \cdot ( \underset{\text{mul}}{D} + \underset{\text{+ op}}{(D-1)} )$

   $= BF(2D-1)$

3. $AI = \dfrac{BF(2D-1)}{BD+DF+BF} \simeq \dfrac{2BDF}{BD+DF+BF}$

   assume $DF$ dominates

   $= 2B$

   if $2B > \dfrac{3.94\times10^{14}}{8.1\times10^{11}}$

4. $T_{\text{math}} = \dfrac{2BDF}{3.94\times10^{14}}$

   $T_{\text{comm}} = \dfrac{BD+DF+BF}{8.1\times10^{11}}$

   $\Leftrightarrow B > 2.43\times10^2$

   $\Leftrightarrow B > 243$

$$\text{Lower Bound} = \max(T_{math}, T_{comm})$$
$$\text{Upper Bound} = T_{math} + T_{comm}$$

(assume we <u>cannot</u> overlap them)

---

## Question 2: int8 + bf16

$$\text{bfloat16}[B, D] \times \text{int8}[D, F] \to \text{bfloat16}[B, F]$$

$$\text{Load} = 2BD + 1 \cdot DF$$

$$\text{Write} = 2BF$$

$$\text{Compute FLOPs} = BF \cdot (\overset{mul}{D} + \overset{+}{(D-1)})$$
$$\simeq 2BDF$$

$$AI = \frac{2BDF}{2BD + DF + 2BF} \simeq \frac{2BDF}{DF} = 2B$$
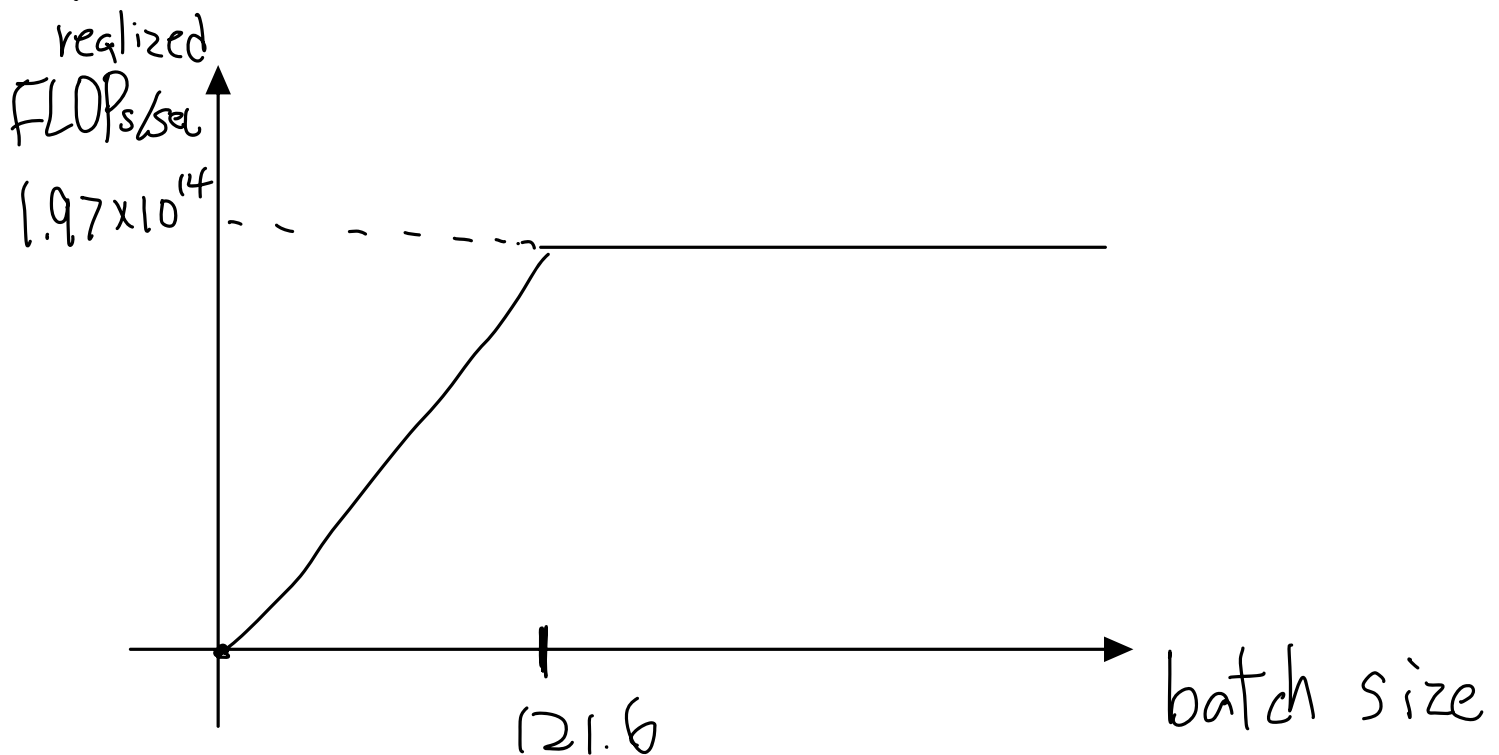
Assume $DF \gg BF$

AI hardware $= \dfrac{1.97 \times 10^{14}}{8.1 \times 10^{11}} = 243.20$    PF>>BD

$\Rightarrow \quad 2B > 243.20$

$\Rightarrow \quad B > 121.6$

Observation: Easier to get to compute-bound when $[D, F]$ matrix is in int8

Question 3:



realized FLOPs/sec

$1.97 \times 10^{14}$

121.6

batch size

(Assume $D, F \gg B$ when $AI \approx \frac{2BDF}{DF}$

$= 2B$

which means $D, F$ do not matter )

If $D, F \gg B$
is not true, we need to consider all
terms in $AI = \frac{2BDF}{2BD + DF + 2BF}$

Pick some random numbers like $D = F = 128$

$AI = \frac{2B \cdot 128 \times 128}{2B \times 128 + 128 \times 128 + 2B \times 128}$

$= \frac{256B}{2B + 128 + 2B} = \frac{256B}{4B + 128} = \frac{64B}{B + 32}$

This means: $\lim_{B \to \infty} \frac{64B}{B+32} = 64 < 243$

$\underset{AI_{hardware}}{\underbrace{}}$

This means when $D = F = 128$, it will always be memory-bound no matter how big is $B$

But when $D = F = 8000$, AI becomes

$$AI = \frac{2B \times 8000 \times 8000}{2B \times 8000 + 8000^2 + 2B \times 8000}$$

$$= \frac{2B \times 8000}{2B + 8000 + 2B}$$

$$= \frac{B \times 8000}{B + 4000 + B}$$

$$= \frac{8000B}{2B + 4000} = \frac{4000B}{B + 2000}$$

When $AI > 243$

$$\iff \frac{4000B}{B + 2000} > 243$$

$\Longleftrightarrow$ $4000B > 243B + 243 \times 2000$

$\Longleftrightarrow$ $3757B > 243 \times 2000$

$\Longleftrightarrow$ $B > 129.35$

This means when $D = F = 8000$,

$\qquad B > 129.35 \Longleftrightarrow$ it is compute-
$\qquad\qquad\qquad\qquad\qquad$ bound

If $D = F = 16,000$

$$\frac{B \times 16000}{B + 8000 + B} = \frac{8000B}{B + 4000} > 243$$
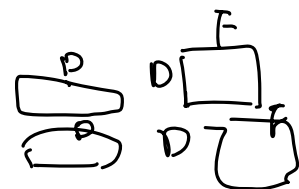
$\Rightarrow$ $8000B > 243B + 243 \times 4000$

$\Rightarrow$ $B > 125.30$

So this matches that if $D, F \gg B$,
$\qquad$ should be converging to $B > 121$

## Question 4:

$$\text{int8}[B, D] *_D \text{int8}[B, D, F]$$
$$\rightarrow \text{int8}[B, F]$$

**Answer:**

There are $B$ $[1, D] \times [D, F]$ matmul

so $\text{FLOPs} = B \cdot (D + (D-1)) \cdot F$
$$\cong 2BDF$$

Load & Write $= BD + BDF + BF$

Arithmetic Intensity $= \dfrac{2BDF}{BD + BDF + BF}$

$$D, F \gg B$$
$$\cong 2$$

This means $\quad 2$ will always $< 243$

so always memory-bound.

---

Question 5:

$$\frac{1.979 \times 10^{15} \text{ FLOPs}}{2} \approx 1 \times 10^{15} \text{ FLOPs/sec}$$

memory bandwith $= 3.35 \times 10^{12}$ bytes/sec

$$AI = \frac{1 \times 10^{15}}{3.35 \times 10^{12}} = \frac{1000}{3.35} = 298$$

$$AI_{\text{bfloat matmul}} = \frac{2BDF}{2DF} = B$$

So $B > 298$ is when it becomes compute-bound