

$$T_{\text{comm}} =$$

$$\frac{\text{Comm bytes}}{\text{bandwidth / sec}}$$

depend on how
you write code

$$T_{\text{math}} =$$

$$\frac{\text{Computation FLOPs}}{\text{Accelerator FLOPs / sec}}$$

Constant,
determined by
hardware

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

(assuming overlap computation)

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

When

$T_{\text{math}} > T_{\text{comm}}$, this means communication
can be overlapped with T_{math} ,
so actual $\frac{\text{FLOPs}}{\text{sec}} = \frac{\text{Computation FLOPs}}{T_{\text{math}}}$

$$= \frac{\text{Accelerator FLOPs}}{\text{sec}}$$

(constant)

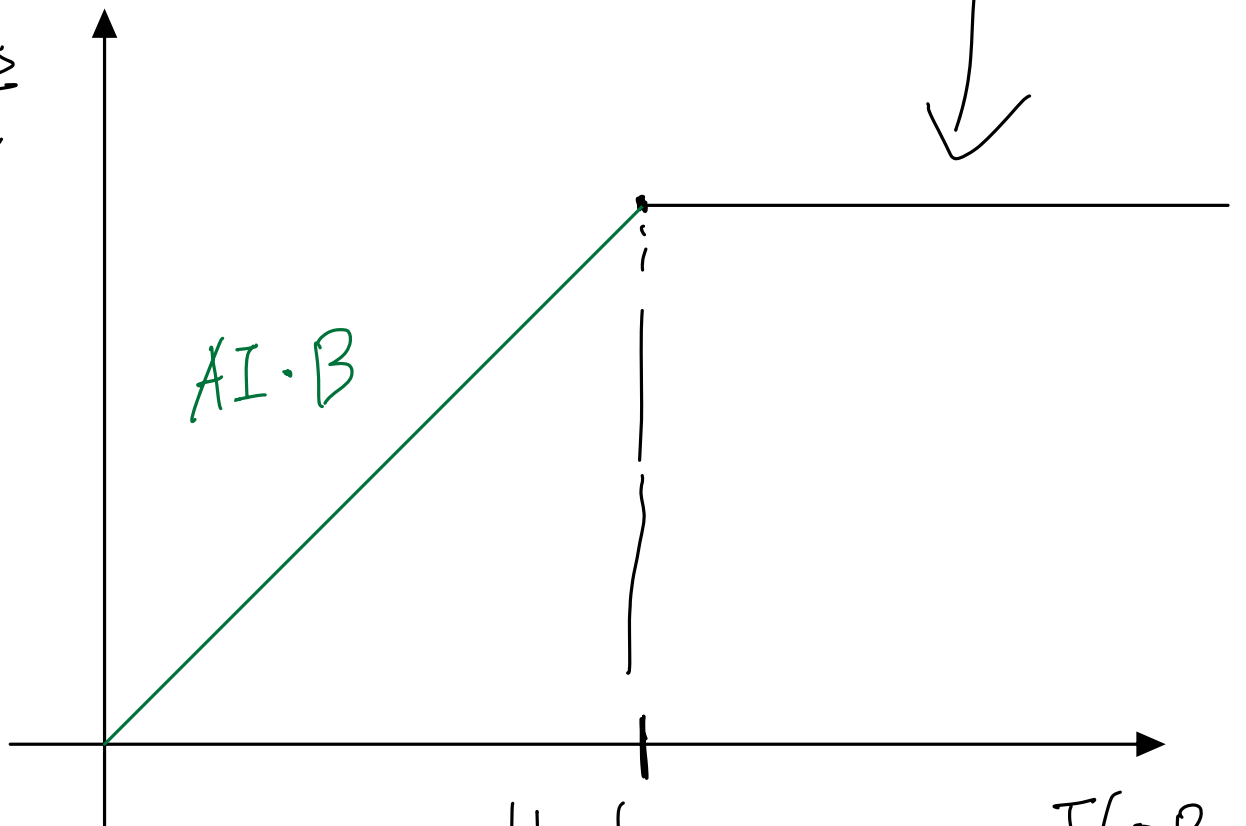
$T_{\text{math}} > T_{\text{comm}}$ also means

$$\Leftrightarrow \frac{\text{Computation FLOPs}}{\text{Accelerator FLOPs}} > \frac{\text{Communication bytes}}{\text{bandwidth bytes/sec}}$$

$$\Leftrightarrow \text{Arithmetic Intensity (compute)} > \text{AI (hardware)}$$

Summary: when $\text{AI}_{\text{compute}} > \text{AI}_{\text{hardware}}$
actual $\frac{\text{FLOPs}}{\text{sec}}$ realized is
always $\text{Accelerator FLOPs/sec}$

Actual
FLOPs
Sec



Hardware
Arithmetic
Intensity

FLOPs
bytes

When

$$T_{\text{math}} < T_{\text{comm}}$$

$$\Leftrightarrow AI_{\text{compute}} < AI_{\text{hardware}}$$

This also means the actual throughput is

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

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

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

This means, as AI_{compute} increase

the realized $\frac{\text{FLOPs}}{\text{Sec}}$ increases linearly

Combining results from:

- ① $T_{\text{math}} < T_{\text{comm}}$
- ② $T_{\text{math}} > T_{\text{comm}}$

observations, we get the
roofline model