

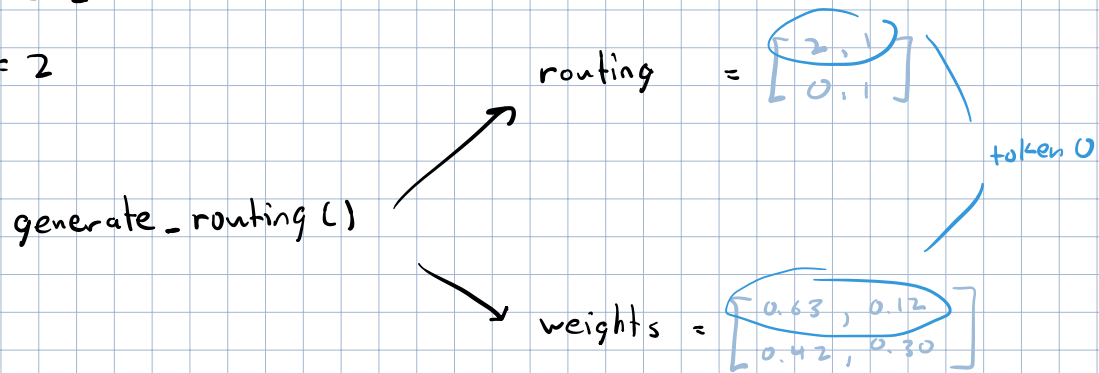
# SIMULATION

## ● HYPERPARAMETERS

✦ Num\_experts = 4

✦ SeqLen = 2

✦ Top-K = 2



✦ Hot\_ratio = 0.5

=> 50% of experts get 80% of routings

✦ Hot\_weight = 0.8

weights  $\rightarrow$  npu\_identify()  $\rightarrow$  labelled weights

$$\begin{bmatrix} 0.63, 0.12 \\ 0.42, 0.30 \end{bmatrix}$$

$$\begin{bmatrix} 0.63, 0.12, (0) (0) \\ 0.42, 0.30, (1) (1) \end{bmatrix}$$

src npu id      token id

## PERFORMANCE MODEL

labelled weights  
+  
routing

→ `convert_to_bytes()` → load dictionary

$\{0: \{0: 0, 1: 2\}, 1: \{0: 3, 1: 0\}\}$

↓        ↓        ↓

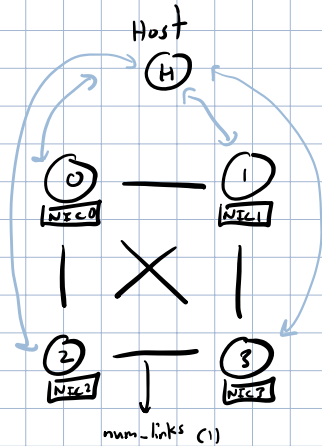
src    dest   load

\* Num\_links = 1

★ Num-nodes = 4

## Mock 1st round

✦ initial cpu delay



loop through load and  
send packets until one or multiple of these:

1) No more load to send (DONE)

2) No more available links to send remaining load

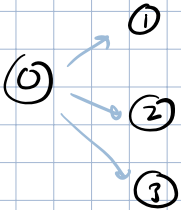
3) No more bandwidth in desired links (leads to 2))

round-time = (most packets in one link) \* packet prep delay

+ (largest bandwidth usage) /  $\rightarrow$  intra-bw

+  $\gamma$  Base\_delay GPU  $\rightarrow$  CPU confirmation

Round robin between dest. nodes for  $\text{round-robin-max-packets per dest}$

$$\begin{aligned} \{0\} &= \{1\} = 3, \\ 2 &= 4, \\ 3 &= 1 \{3\} \end{aligned}$$


#1) 0-1 2B

42) 0-2 2B

# 31003 1B

#4) 0 ~ 1 1B

#5) 0-2 2B

$$\frac{11}{2}$$

Assuming a balanced routing,

$$\text{Round-time} = (\text{initial-cpu-delay} + \text{base-delay} + \left\lfloor \frac{\text{intra-bw}}{\text{packet-size}} \right\rfloor) \left( \text{packet-prep-delay} + \frac{\text{packet-size}}{\text{intra-bw}} \right)$$

$$\text{Num-round} = \frac{\text{load\_per\_node} \cdot \text{num\_nodes}}{\left( \frac{\text{num\_links} \cdot \text{num\_nodes} \cdot (\text{num\_nodes} - 1)}{2} \right) \cdot \text{intra-bw}}$$

$$\text{Total\_time} = \text{Round-time} \cdot \text{Num-rounds}$$

$$\begin{aligned} \text{load} &= \{0: \{0:0, 1:10, 2:10\}, \\ &\quad 1: \{0:10, 1:0, 2:10\}, \\ &\quad 2: \{0:10, 1:10, 2:0\}\} \end{aligned}$$

↓

$$\text{load\_per\_node} = 20$$