# SIMULATION

⇒ Num_experts = 4

⇒ Seqlen = 2

⇒ Top_k = 2

generate_routing ()

$$\text{routing} = \begin{bmatrix} 2, 1 \\ 0, 1 \end{bmatrix}$$

token 0

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

⇒ Hot_ratio = 0.5

⇒ Hot_weight = 0.8

⇒ 50% of experts get 80% of routings

weights ⟿ npu_identify () ⟿ 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
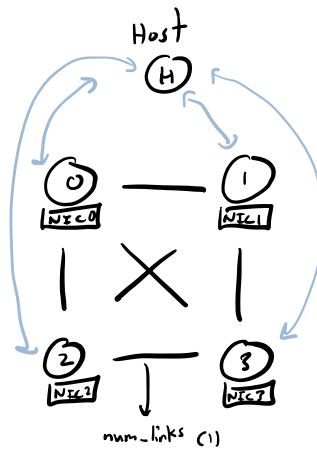$\longrightarrow$ convert_to_bytes() $\longrightarrow$ load dictionary

$\{0: \{0:0, 1:2\}, 1: \{0:3, 1:0\}\}$
      ↓      ↓    ↓
     src   dest  load

✷ Num_links = 1

✷ Num_nodes = 4

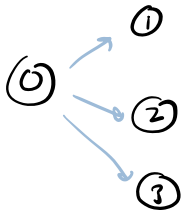## Mock 1st round

Host



✷ initial cpu delay {

num_links (1)

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) / ✷ intra_bw

+ ✷ Base_delay    GPU ⟶ CPU confirmation

## Load

$\{0: \{1:3, 2:4, 3:1\}\}$



## Order of packets sent

#1) 0→1  2B

#2) 0→2  2B

#3) 0→3  1B

#4) 0→1  1B

#5) 0→2  2B

Assuming a balanced routing,

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

$$\text{Num-round} = \frac{\ast\,\text{load-per-node} \cdot \ast\,\text{num-nodes}}{\left( \frac{\ast\,\text{num-links} \cdot \ast\,\text{num-nodes} \cdot (\ast\,\text{num-nodes} - 1)}{2} \right) \cdot \ast\,\text{intra-bw}}$$

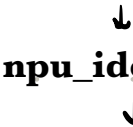$$\text{Total-time} = \text{Round-time} \cdot \text{Num-rounds}$$

$\ast$ load = { 0: { 0:0, 1:10, 2:10 },
                1: { 0:10, 1:0, 2:10 },
                2: { 0:10, 1:10, 2:0 } }

⇓

load_per_node = 20

## Simulation

- NUM_EXPERTS

- TOP_K

### Imbalance

- HOT_RATIO

- HOT_WEIGHT

**generate_routing()**

↓

**npu_identify()** ← identifies by token & src node

↓

## Performance Model

### 1) Load conversion

- Data conversion parameters

  **convert_to_bytes()**

→

### 2) Communication

- Infrastructure parameters

  **full_mesh_comm()**

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

↓

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

src npu id        token id