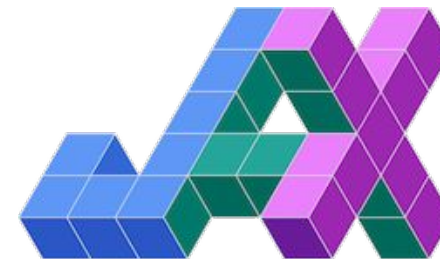


Tensor Parallelism in



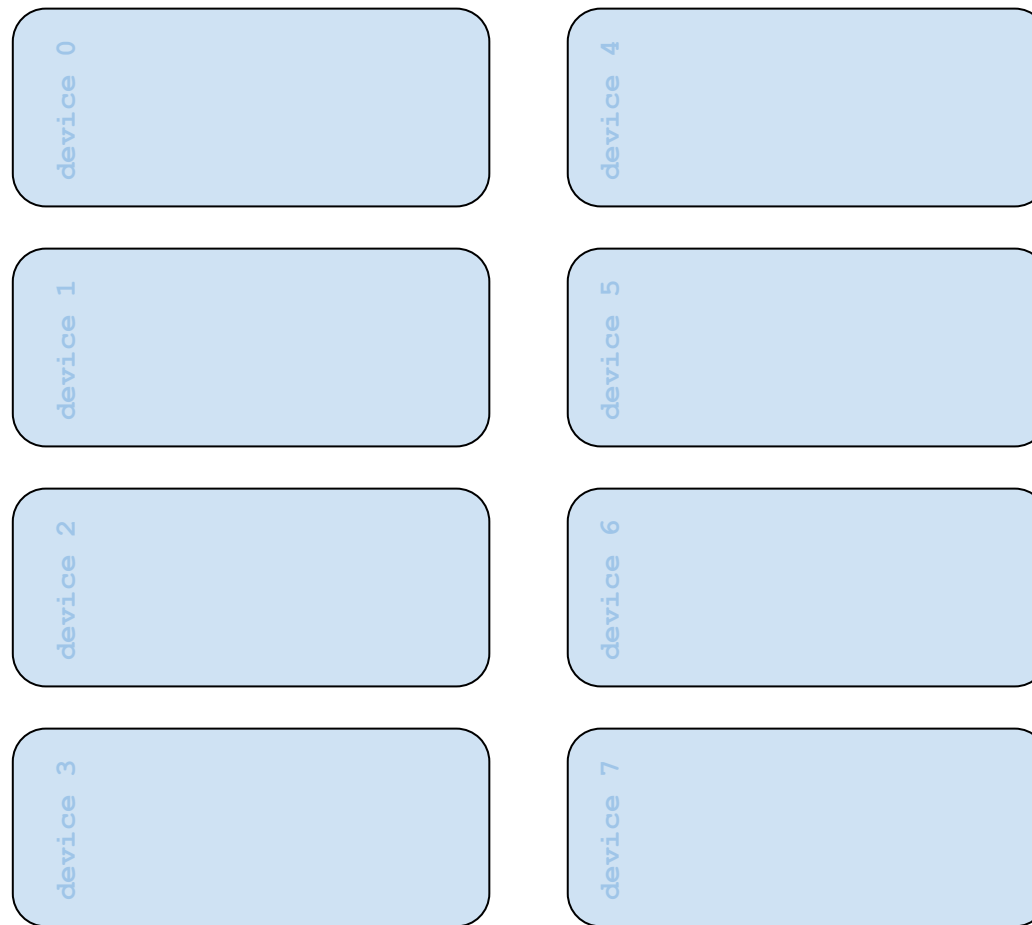
Eastern European Machine Learning Summer School
10-15 July 2023, Košice, Slovakia

Outline of tutorial

- Types of parallelism
- Refresher on matrix multiplication
- Megatron-like sharded MLP
- 2D parallelism

Types of parallelism

- Let's assume we have access to 8 interconnected devices that we can arrange in a grid.
- For example, the grid can be (8, 1) or (4, 2)



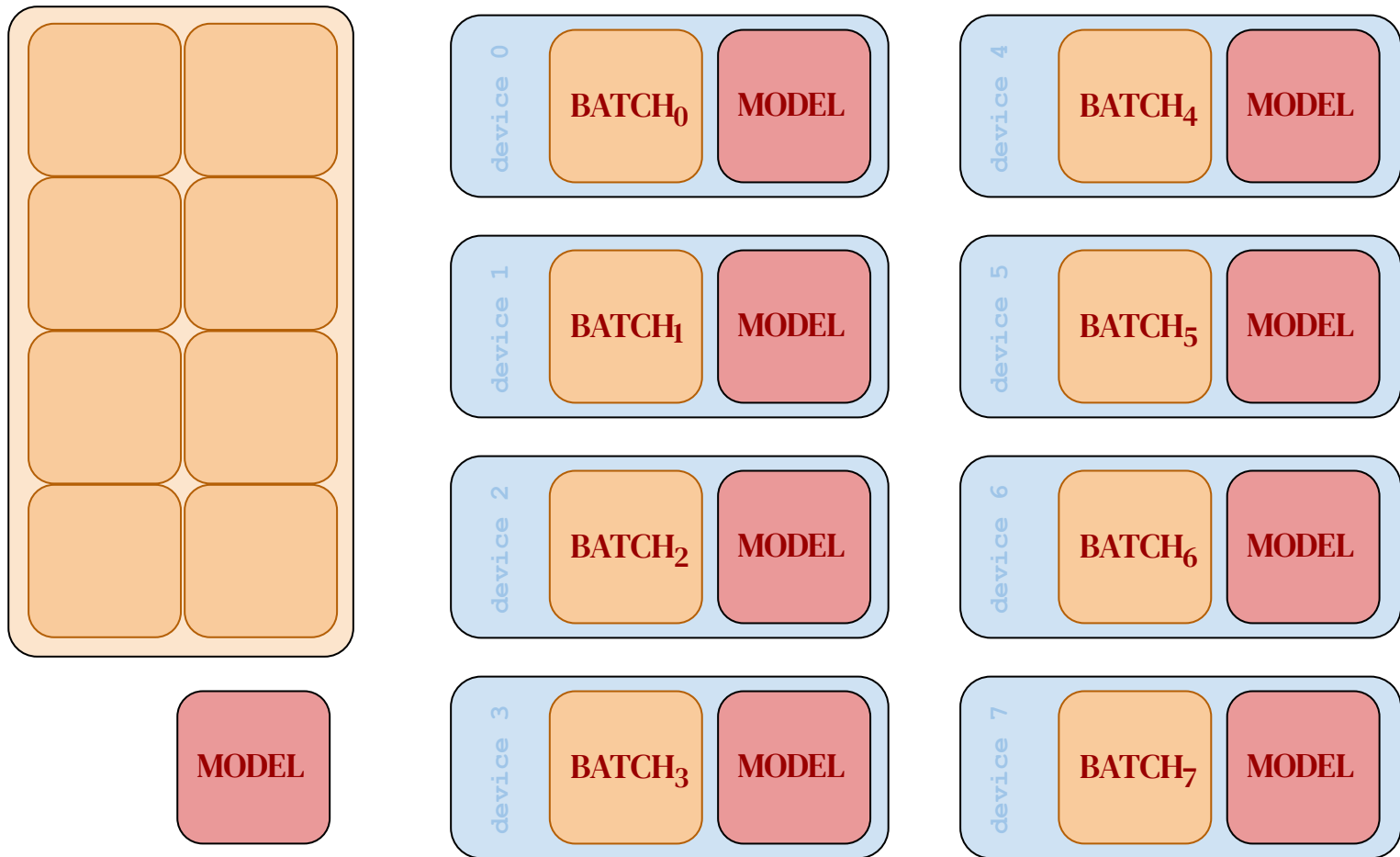
Data parallelism

In this paradigm, we have a model that fits in one device, but we want to speed up training by using larger batches.

To do so, we replicate model across every device and chunk global batch into smaller units.

The same model (forward and backward pass) is executed on each devices, but to keep params in sync, we need to do synchronization step before update.

LARGE BATCH



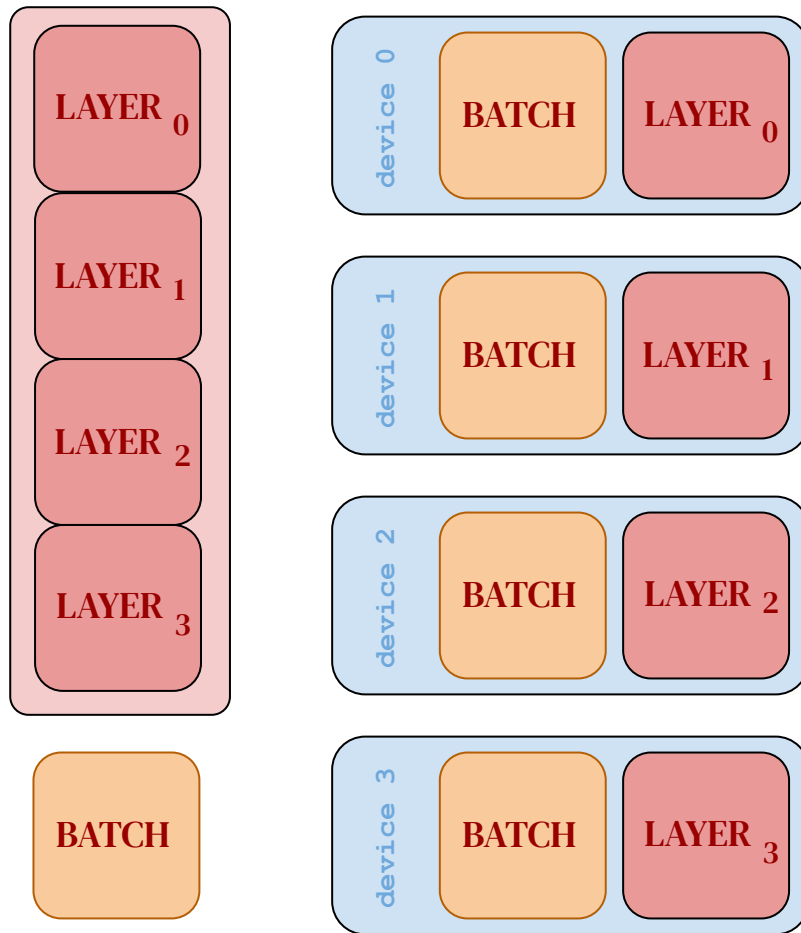
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



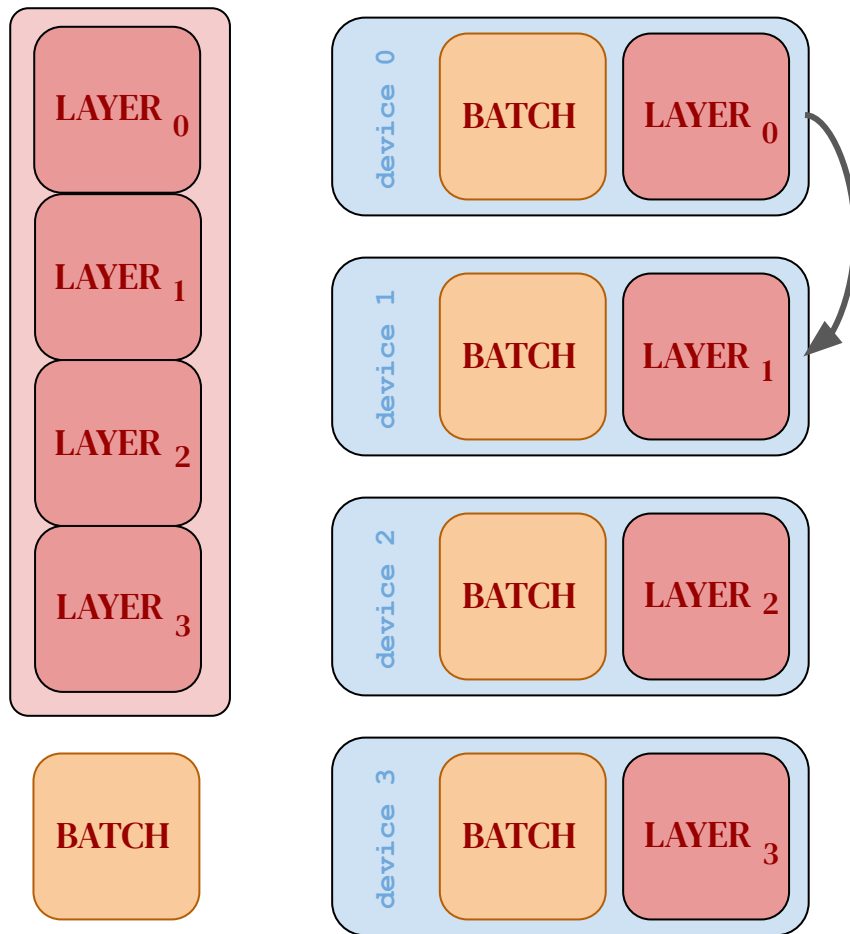
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



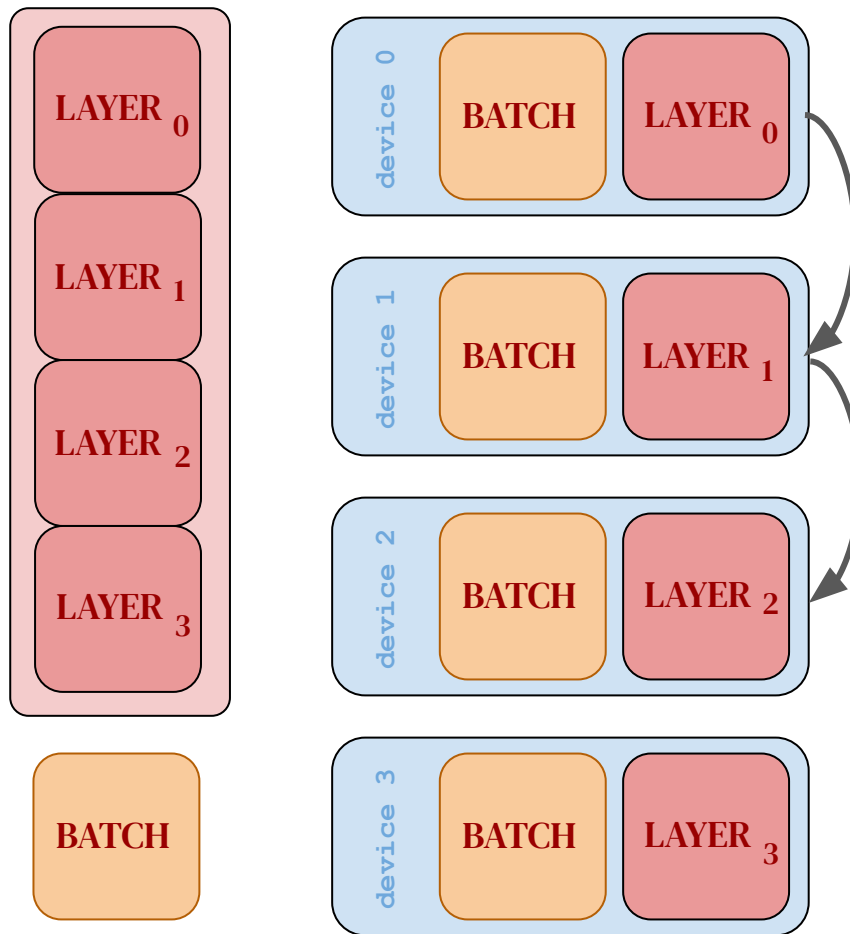
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



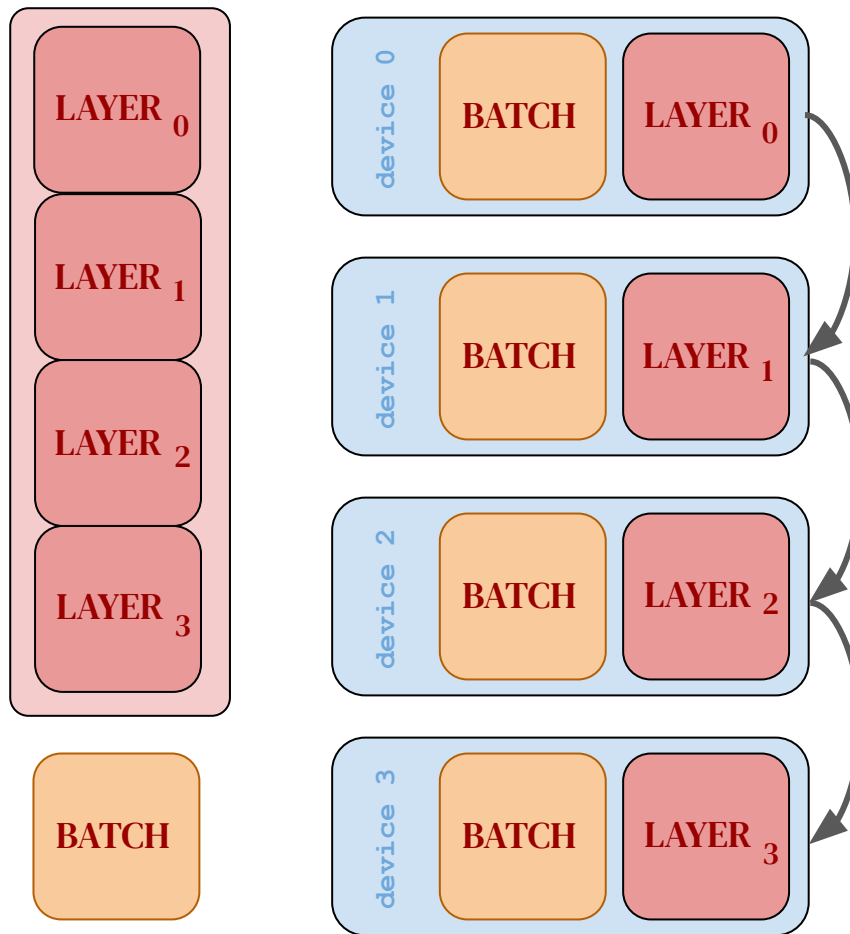
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



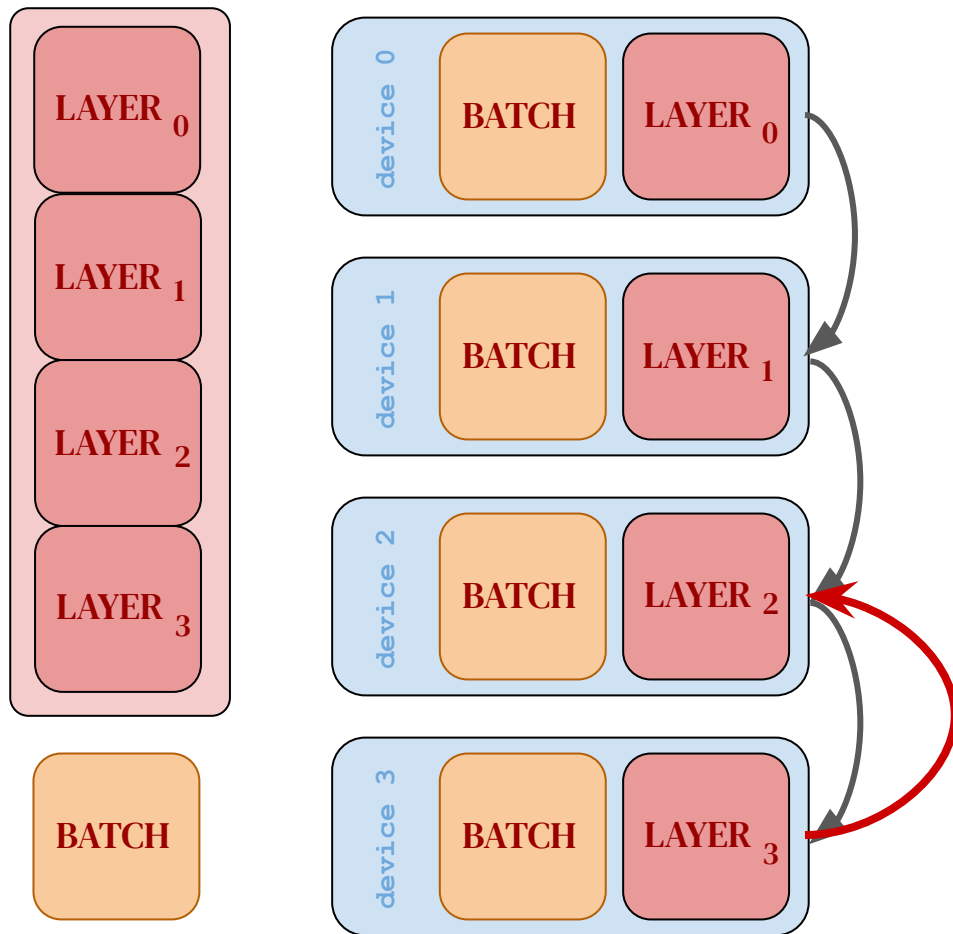
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



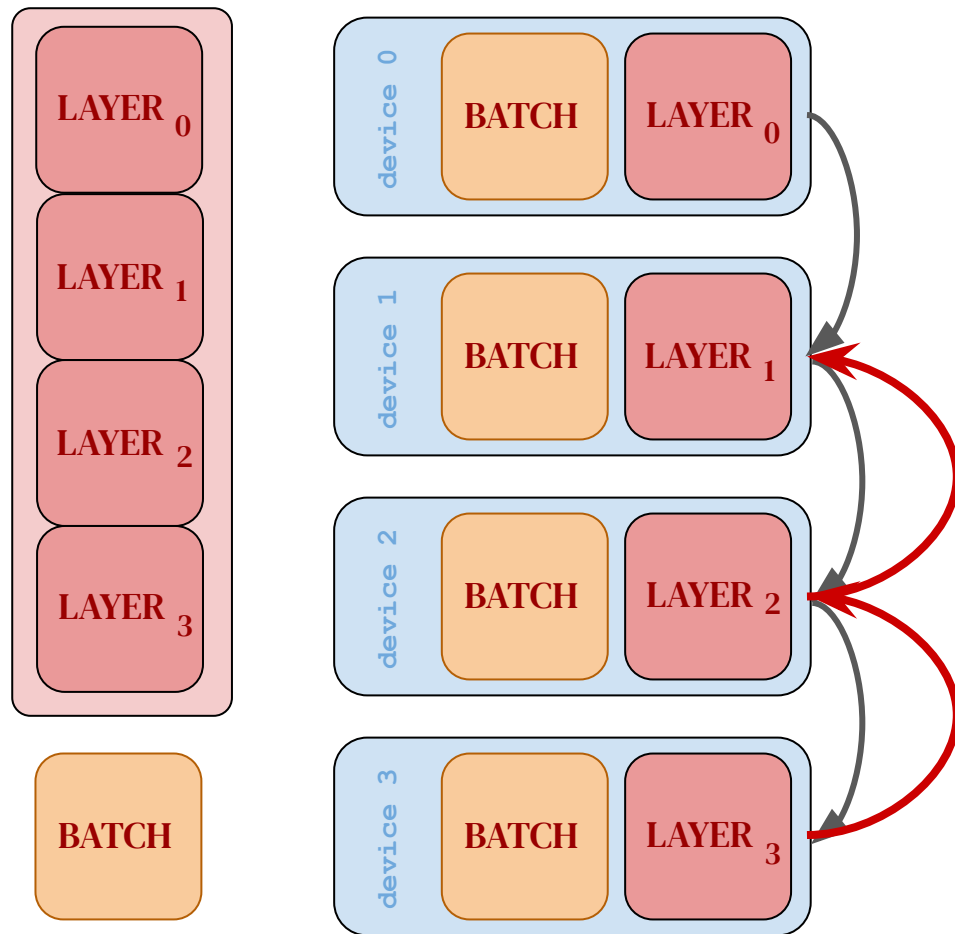
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



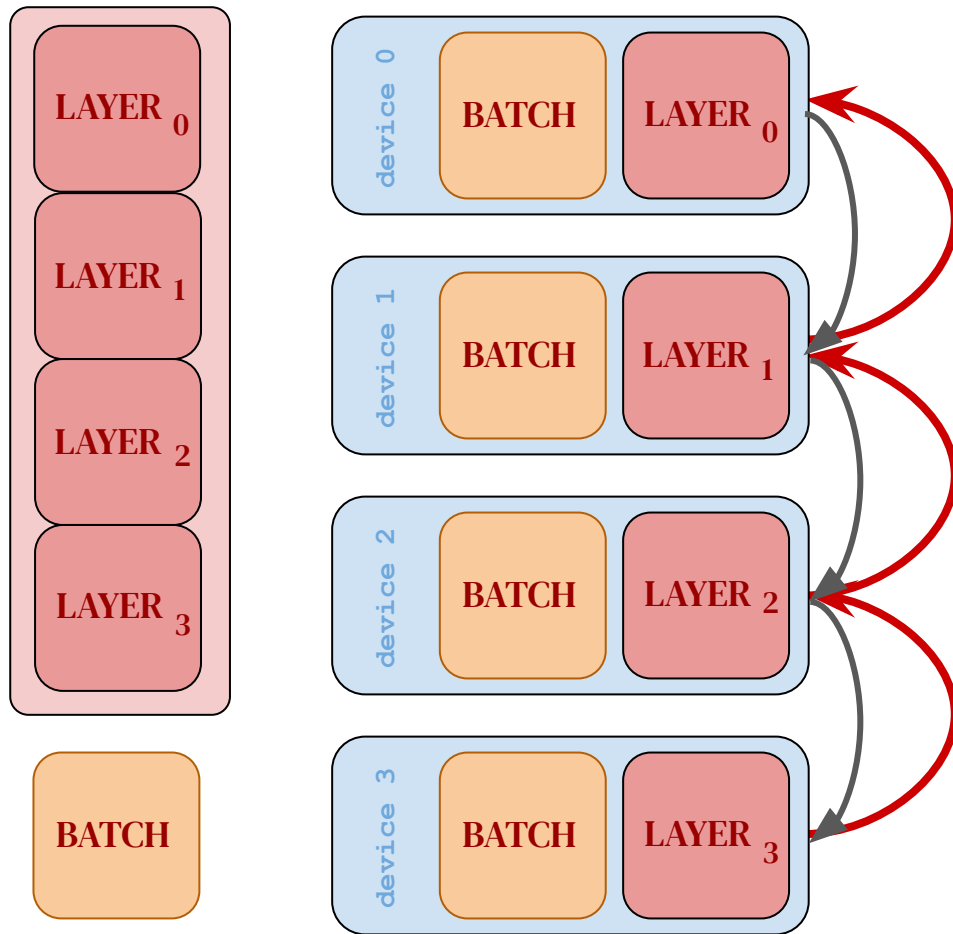
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Pipeline parallelism:** different sub-sequences of layers live on separate accelerators, and the computation flows sequentially from device to device (example [GPipe](#))
- The main difficulty is keeping devices busy.

LARGE MODEL



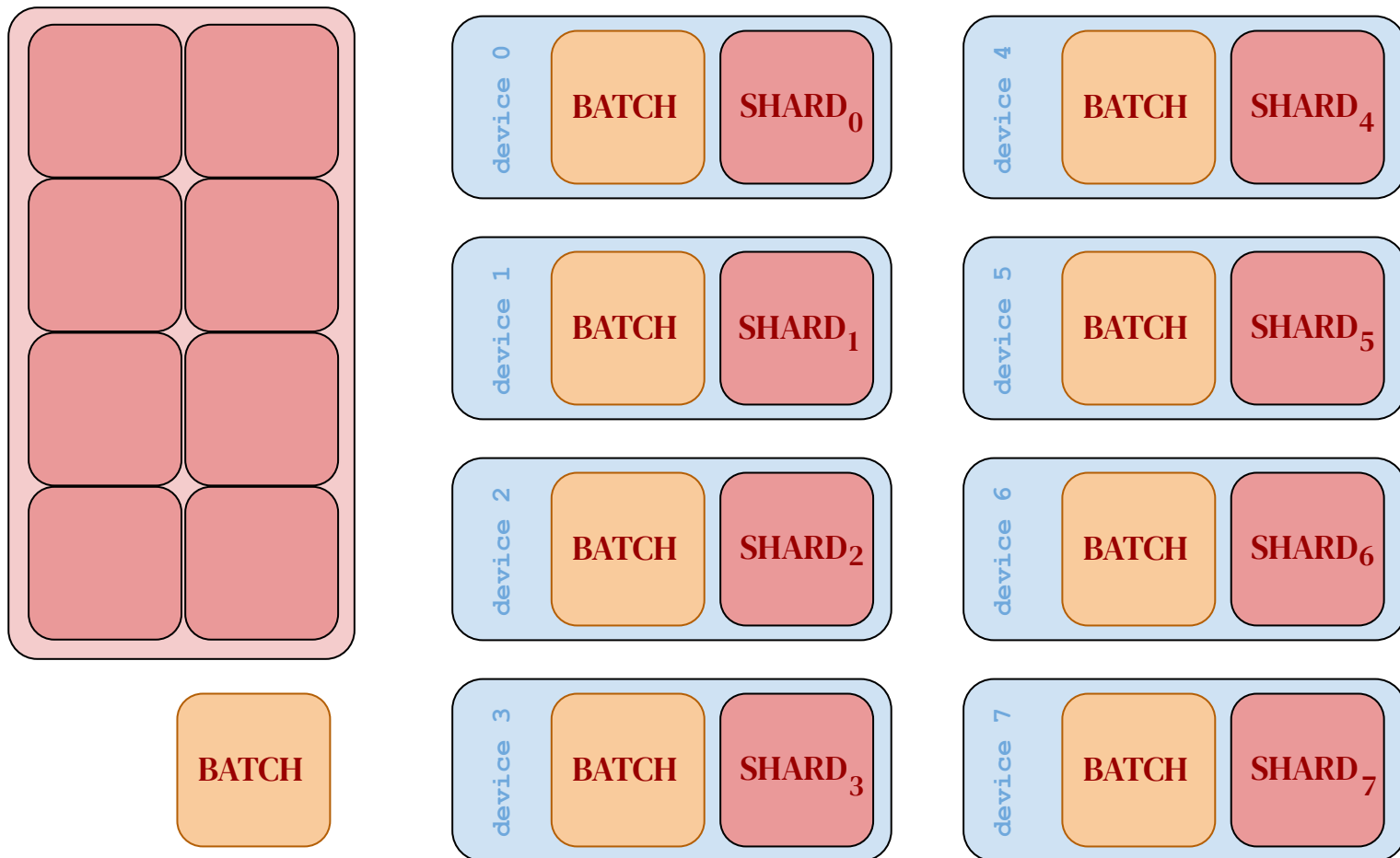
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Tensor parallelism:** individual tensors (activations/weights) are split across devices (example: [GSPMD](#))
- The main difficulty is keeping the communication volume down.

LARGE MODEL



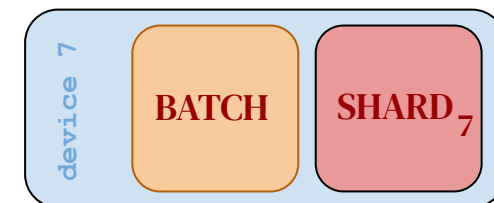
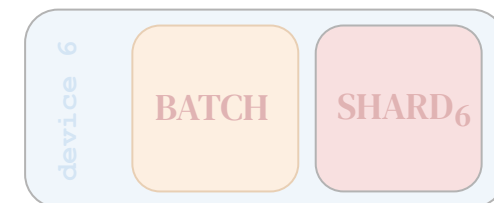
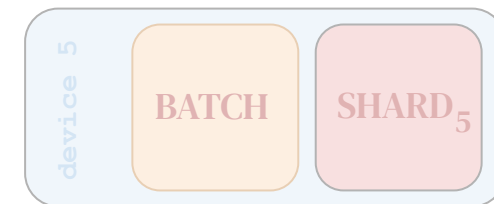
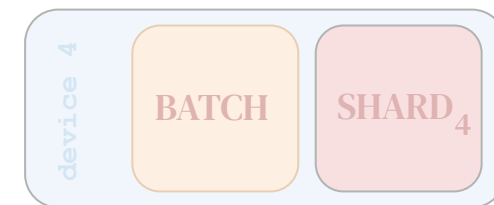
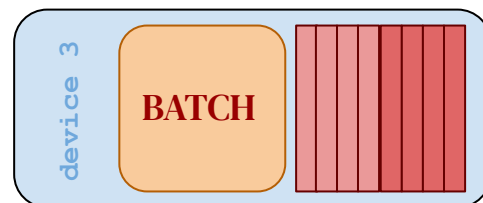
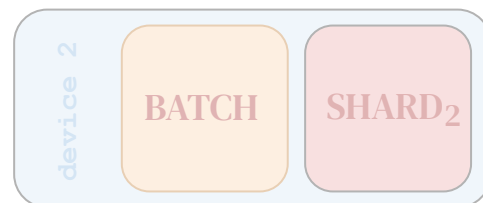
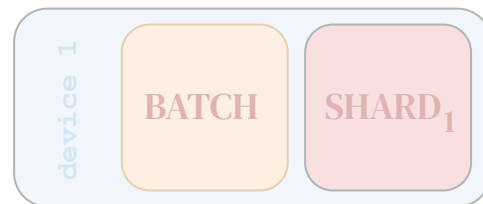
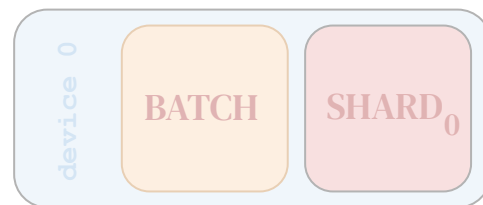
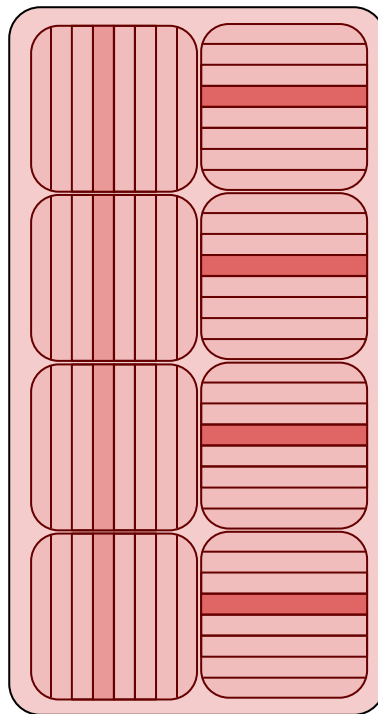
Model parallelism

If our model doesn't fit on one device, we can split computation across multiple devices.

There are different kinds of model parallelism

- **Tensor parallelism:** individual tensors (activations/weights) are split across devices (example: [GSPMD](#))
- The main difficulty is keeping the communication volume down.

LARGE MODEL

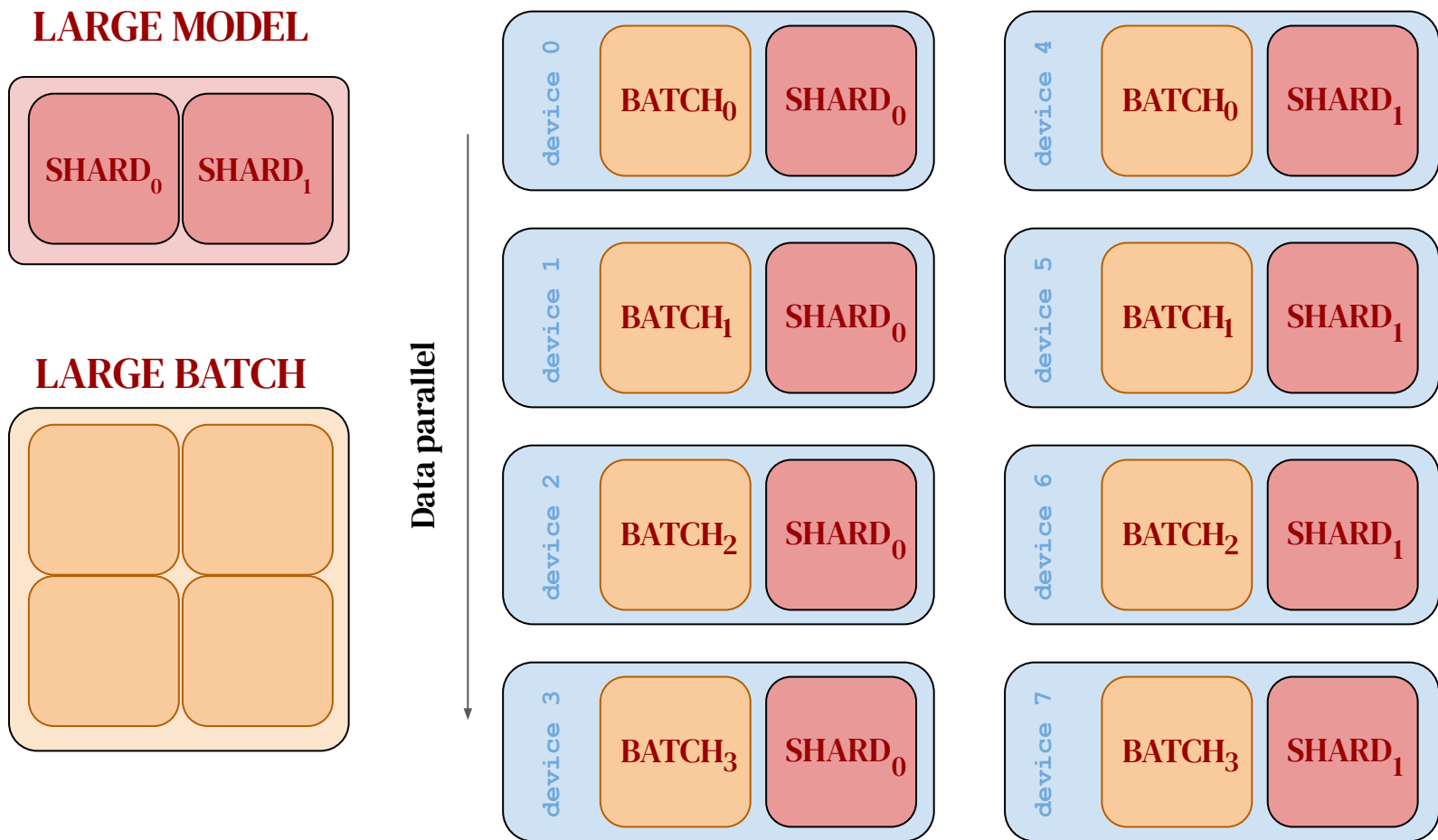


Batch and data parallelism

We can combine data and model parallelism, by arranging our devices in a 2D grid, where each axis represents different partitioning.

For example given (4, 2) grid:

- On the vertical axis we apply **batch parallelism** by splitting our large batch in 4
- On the horizontal axis we apply **model parallelism** by splitting the model in 2.



A Refresher on Matrix Multiplication

The diagram illustrates the calculation of the element C_{21} in matrix C through the dot product of the second row of matrix A and the first column of matrix B .

Matrix A (2x3):

A_{11}	A_{12}	A_{13}
A_{21}	A_{22}	A_{23}

Matrix B (3x1):

B_{11}
B_{21}
B_{31}

Calculation of C_{21} :

$$\begin{aligned} & \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \end{bmatrix} \times \begin{bmatrix} B_{11} \\ B_{21} \\ B_{31} \end{bmatrix} = \begin{bmatrix} A_{11} \times B_{11} + A_{12} \times B_{21} + A_{13} \times B_{31} \\ A_{21} \times B_{11} + A_{22} \times B_{21} + A_{23} \times B_{31} \end{bmatrix} = \begin{bmatrix} C_{11} \\ C_{21} \end{bmatrix} \end{aligned}$$

The final result C_{21} is highlighted in an orange box.

A Refresher on Matrix Multiplication

The diagram illustrates the calculation of the element C_{21} in matrix multiplication. It shows the following steps:

- Matrix A (2x3) is multiplied by Matrix B (3x1).
- The result is the sum of three products: $A_{11} \times B_{11} + A_{12} \times B_{21} + A_{13} \times B_{31}$ (top row) and $A_{21} \times B_{11} + A_{22} \times B_{21} + A_{23} \times B_{31}$ (bottom row).
- The final result is Matrix C (2x1), where the bottom element is C_{21} .

The diagram uses red boxes for the input matrices and intermediate products, and orange boxes for the final result matrix C .

A Refresher on Matrix Multiplication

The diagram illustrates the calculation of the element C_{21} in matrix C through the dot product of the second row of matrix A and the first column of matrix B .

Matrix A (3x3):

A_{11}	A_{12}	A_{13}
A_{21}	A_{22}	A_{23}

Matrix B (3x1):

B_{11}
B_{21}
B_{31}

Calculation of C_{21} :

The second row of A is multiplied by the first column of B using the distributive property:

$$A_{21} \times B_{11} + A_{22} \times B_{21} + A_{23} \times B_{31}$$

The result of this calculation is the element C_{21} in the first row of matrix C .

Matrix C (1x2):

C_{11}
C_{21}

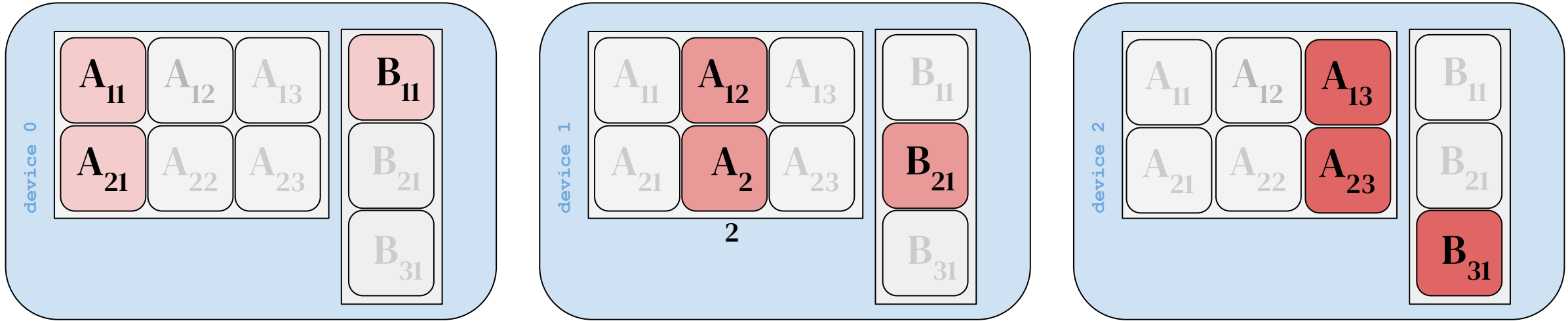
A Refresher on Matrix Multiplication

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \end{bmatrix} \times \begin{bmatrix} B_{11} \\ B_{21} \\ B_{31} \end{bmatrix} = \begin{bmatrix} A_{11} \times B_{11} + A_{12} \times B_{21} + A_{13} \times B_{31} \\ A_{21} \times B_{11} + A_{22} \times B_{21} + A_{23} \times B_{31} \end{bmatrix} = \begin{bmatrix} C_{11} \\ C_{21} \end{bmatrix}$$

Those products can be calculated on separate devices

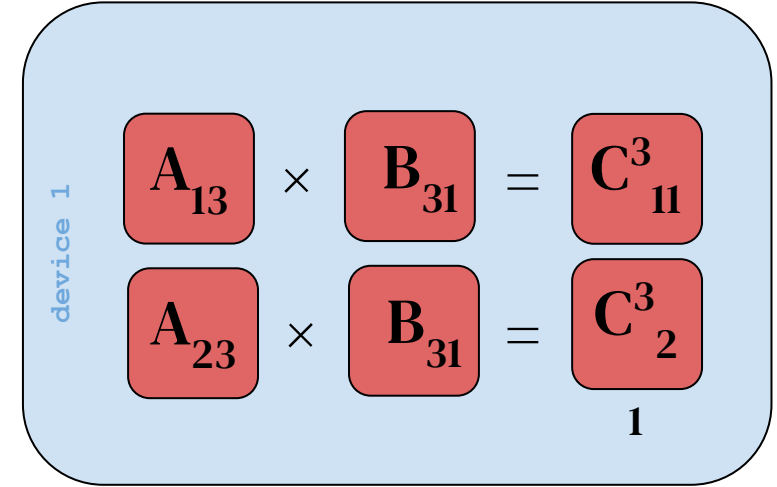
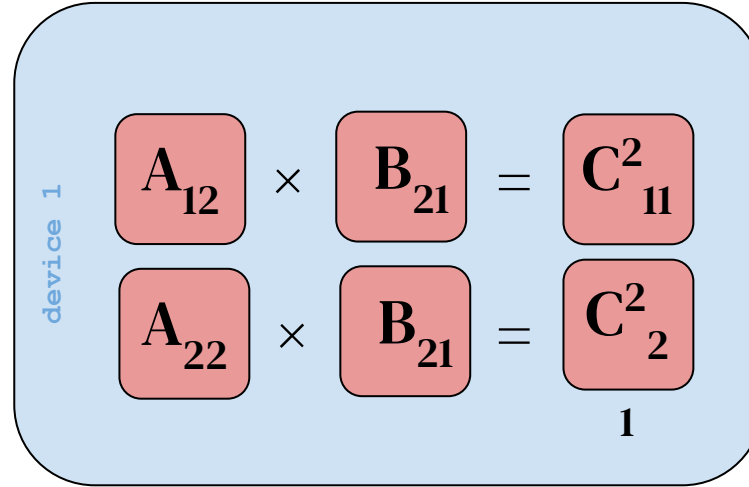
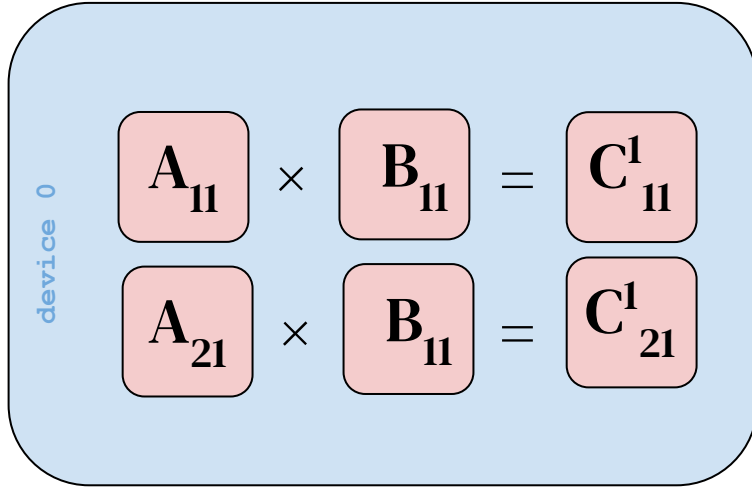
Distributed Matrix Multiplication

Step 1: Shard A on columns, B on rows



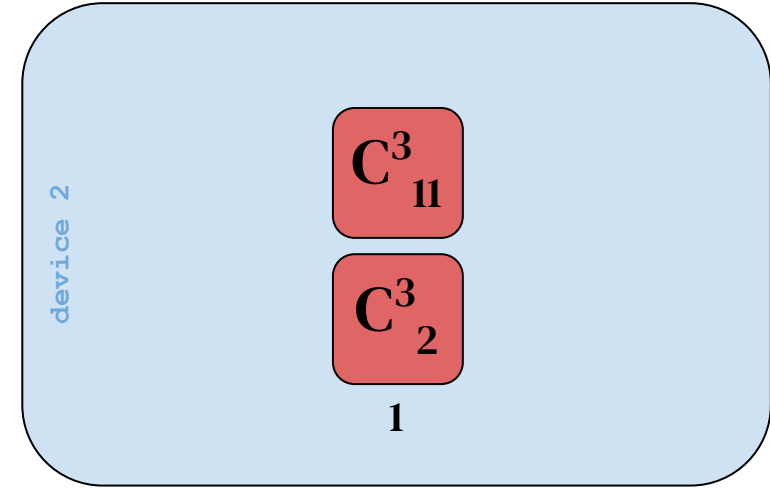
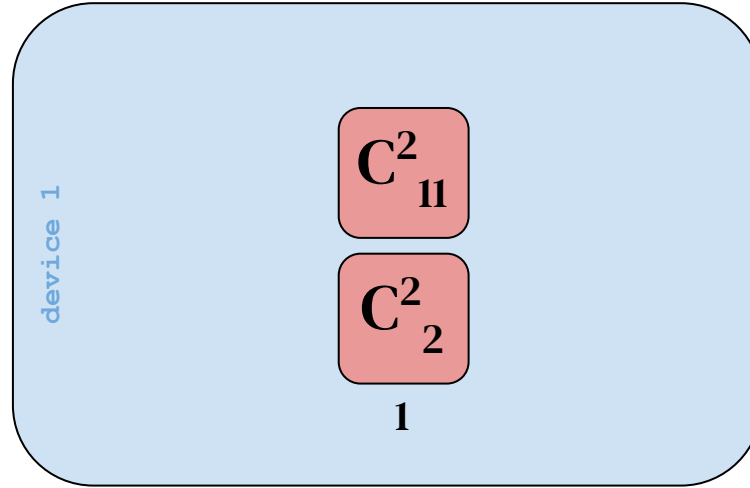
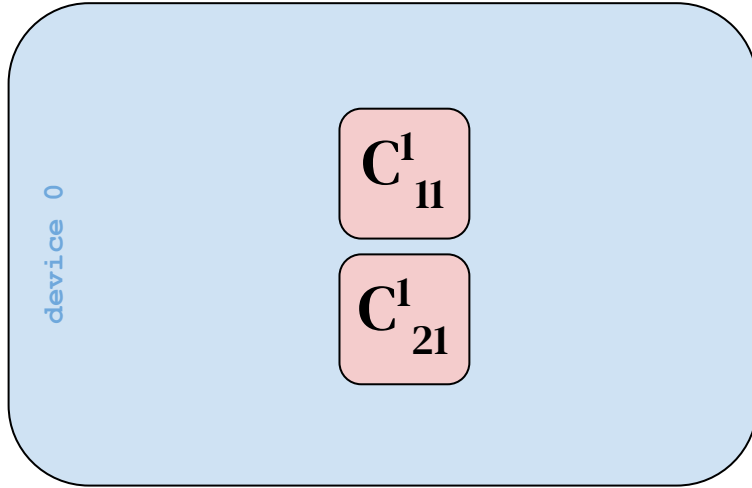
Distributed Matrix Multiplication

Step 2: Perform partial products on devices



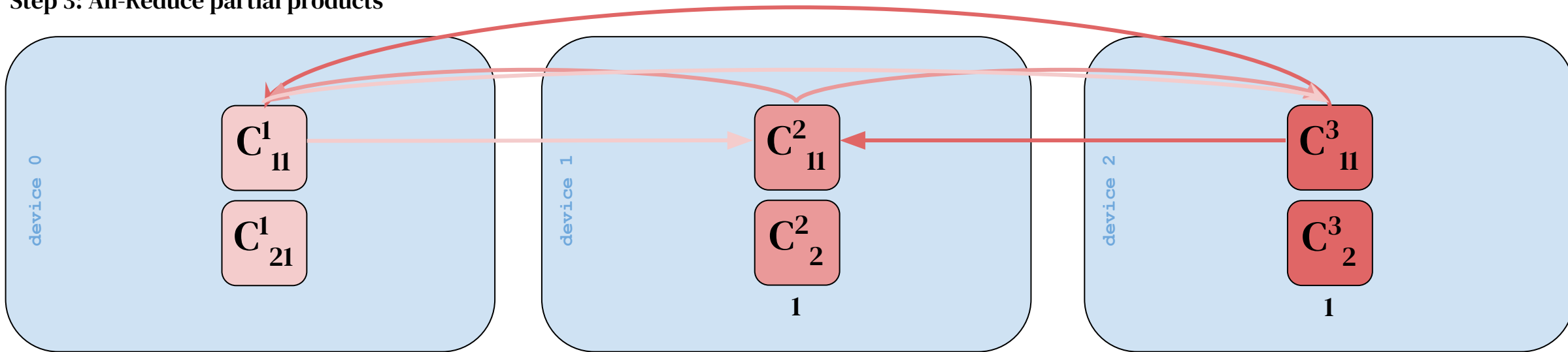
Distributed Matrix Multiplication

Step 3: All-Reduce partial products



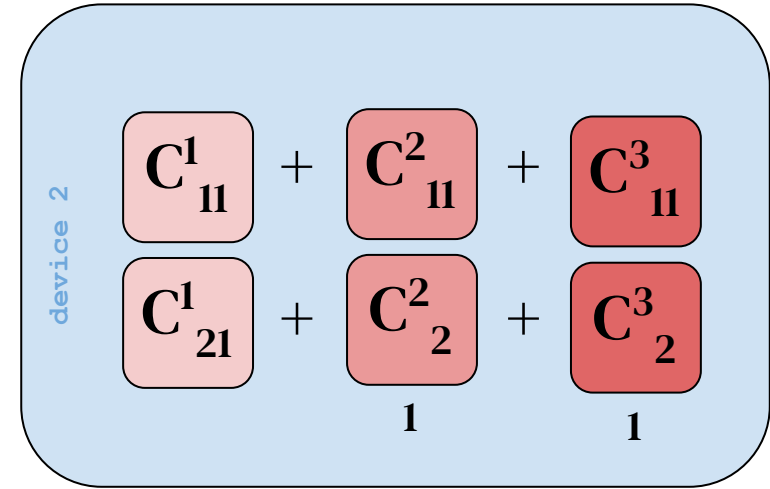
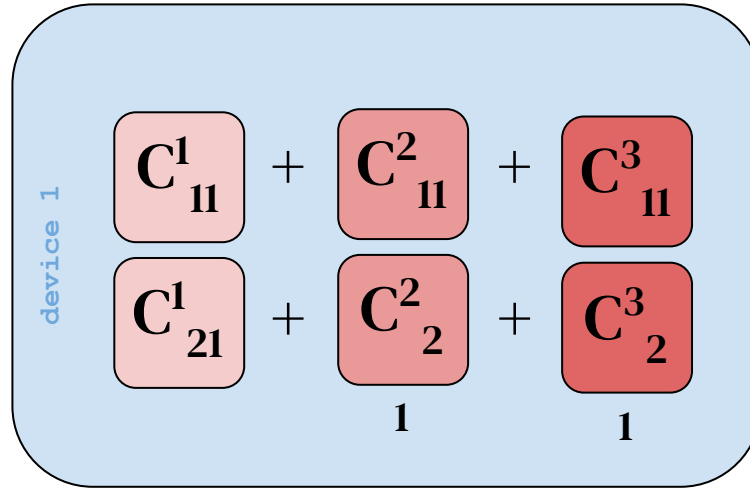
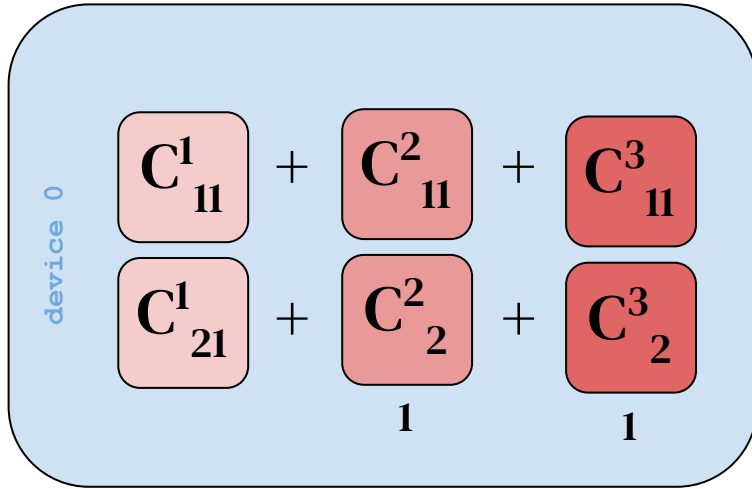
Distributed Matrix Multiplication

Step 3: All-Reduce partial products



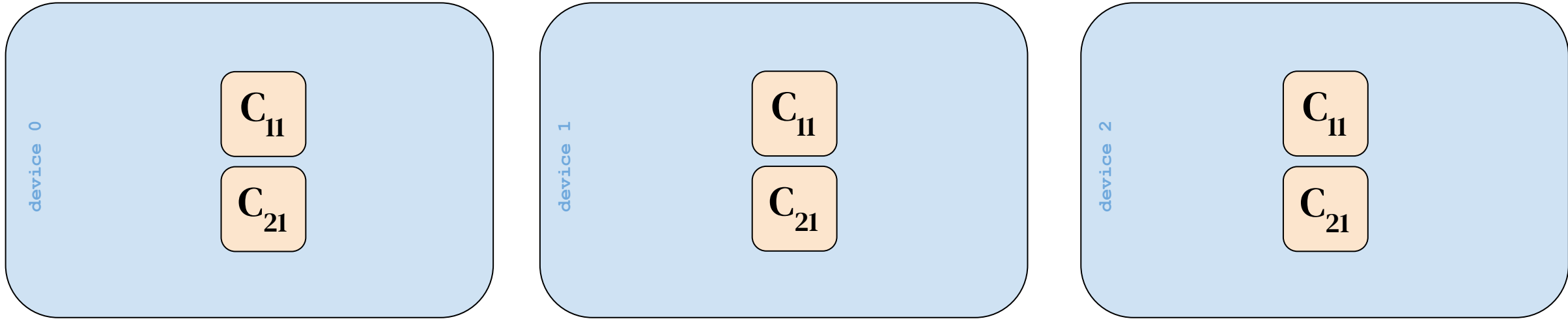
Distributed Matrix Multiplication

Step 3: All-Reduce partial products



Distributed Matrix Multiplication

Step 3: All-Reduce partial products



Distributed Matrix Multiplication - Summary

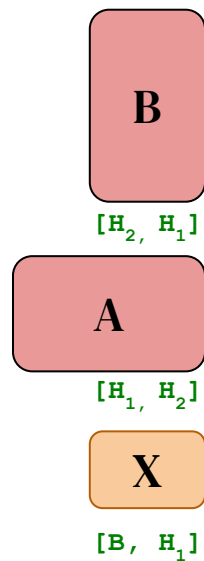
- When multiplying two matrices \mathbf{A} , \mathbf{B} such that their reduction dimension is nicely divisible by number of devices we can split the computation and perform the partial products \mathbf{C}_i in parallel.
- To recover the unsharded $\mathbf{C} = \mathbf{AB}$ we need to all-reduce partial products incurring communication cost.
- This is a recurring pattern in tensor parallelism: we're trading device memory utilization for communication cost.

Sharding a 2-layer MLP

```
def mlp(  
    x: Float[Array, "B H_1"],  
    A: Float[Array, "H_1 H_2"],  
    B: Float[Array, "H_2 H_1"]) -> Float[Array, "B H_1"]:  
    """ Z = max(X * W_1, 0) * W_2 """  
    y = jnp.dot(x, A)      # [B, H_1] @ [H_1, H_2] -> [B, H_2]  
    u = jnp.maximum(y, 0)  # [B, H_2]  
    z = jnp.dot(u, B)      # [B, H_2] @ [H_2, H_1] -> [B, H_1]  
    return z
```

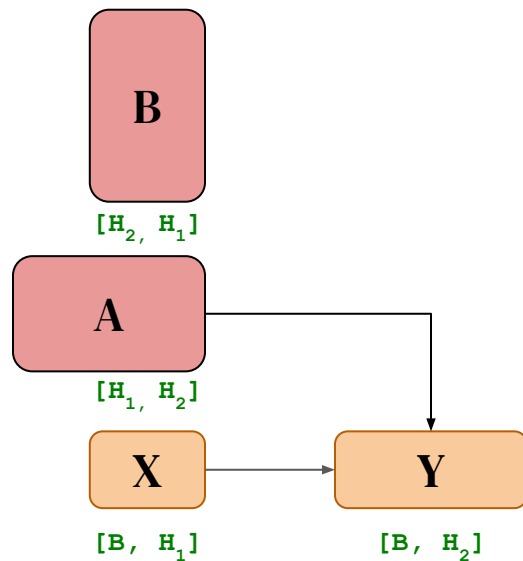
Sharding a 2-layer MLP

```
def mlp(  
    x: Float[Array, "B H_1"],  
    A: Float[Array, "H_1 H_2"],  
    B: Float[Array, "H_2 H_1"]) -> Float[Array, "B H_1"]:  
    """ Z = max(X * W_1, 0) * W_2 """  
    y = jnp.dot(x, A)      # [B, H_1] @ [H_1, H_2] -> [B, H_2]  
    u = jnp.maximum(y, 0)  # [B, H_2]  
    z = jnp.dot(u, B)      # [B, H_2] @ [H_2, H_1] -> [B, H_1]  
    return z
```



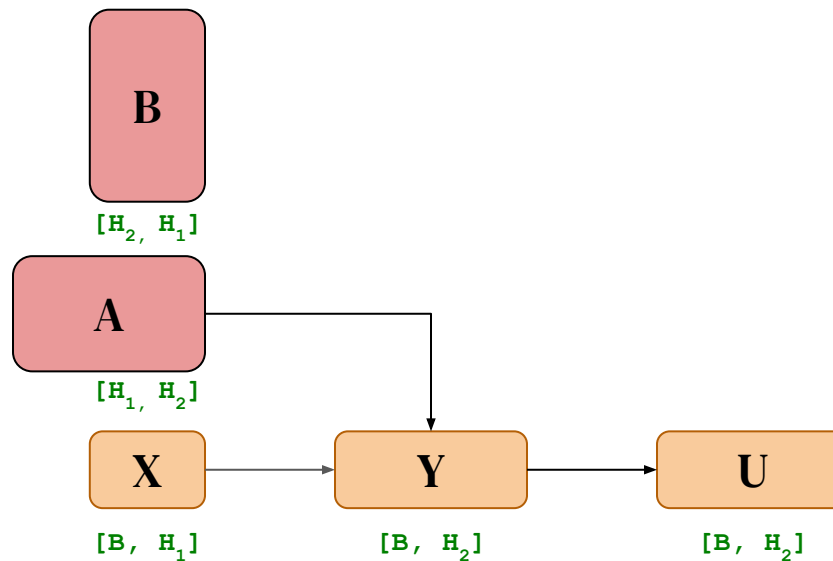
Sharding a 2-layer MLP

```
def mlp(  
    x: Float[Array, "B H_1"],  
    A: Float[Array, "H_1 H_2"],  
    B: Float[Array, "H_2 H_1"]) -> Float[Array, "B H_1"]:  
    """ Z = max(X * W_1, 0) * W_2 """  
    y = jnp.dot(x, A)      # [B, H_1] @ [H_1, H_2] -> [B, H_2]  
    u = jnp.maximum(y, 0)  # [B, H_2]  
    z = jnp.dot(u, B)      # [B, H_2] @ [H_2, H_1] -> [B, H_1]  
    return z
```



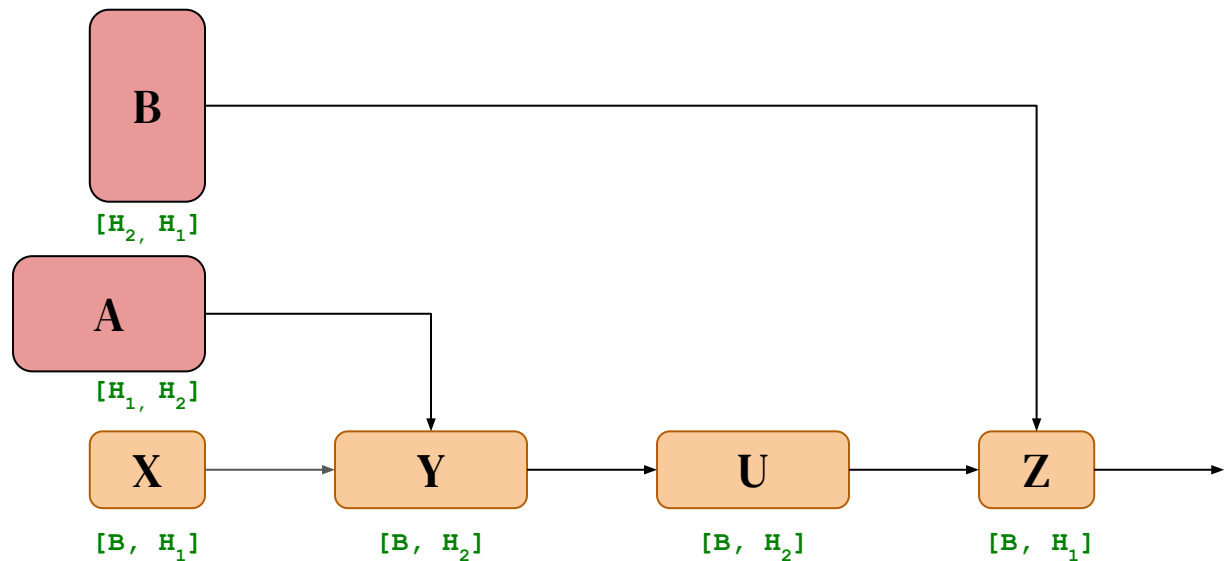
Sharding a 2-layer MLP

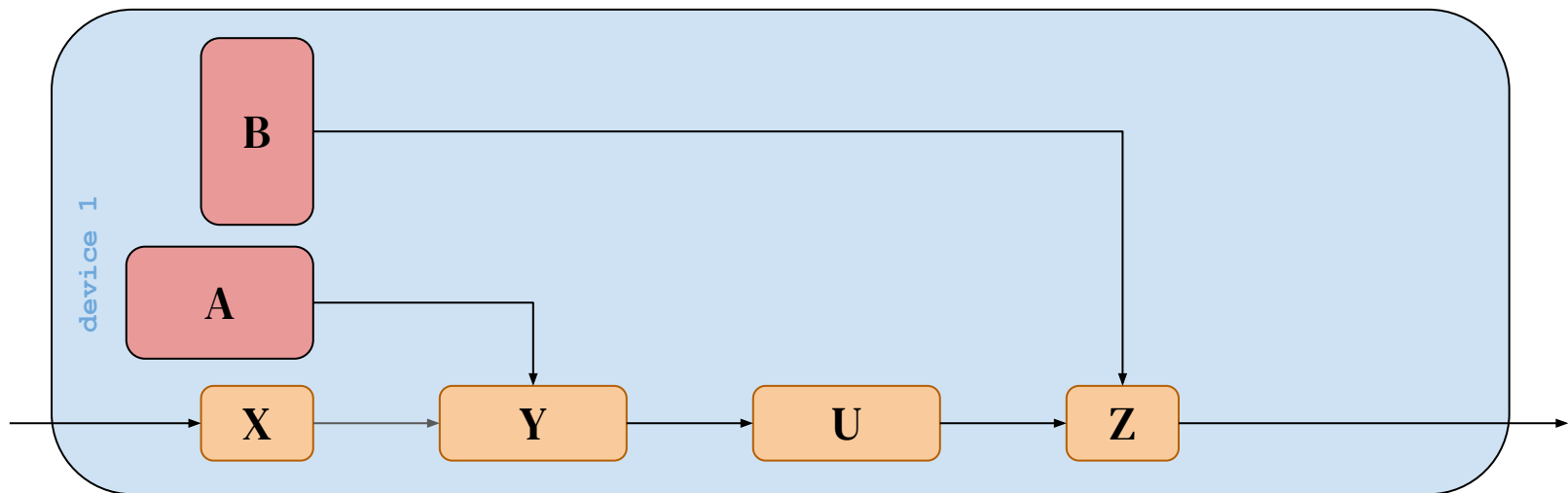
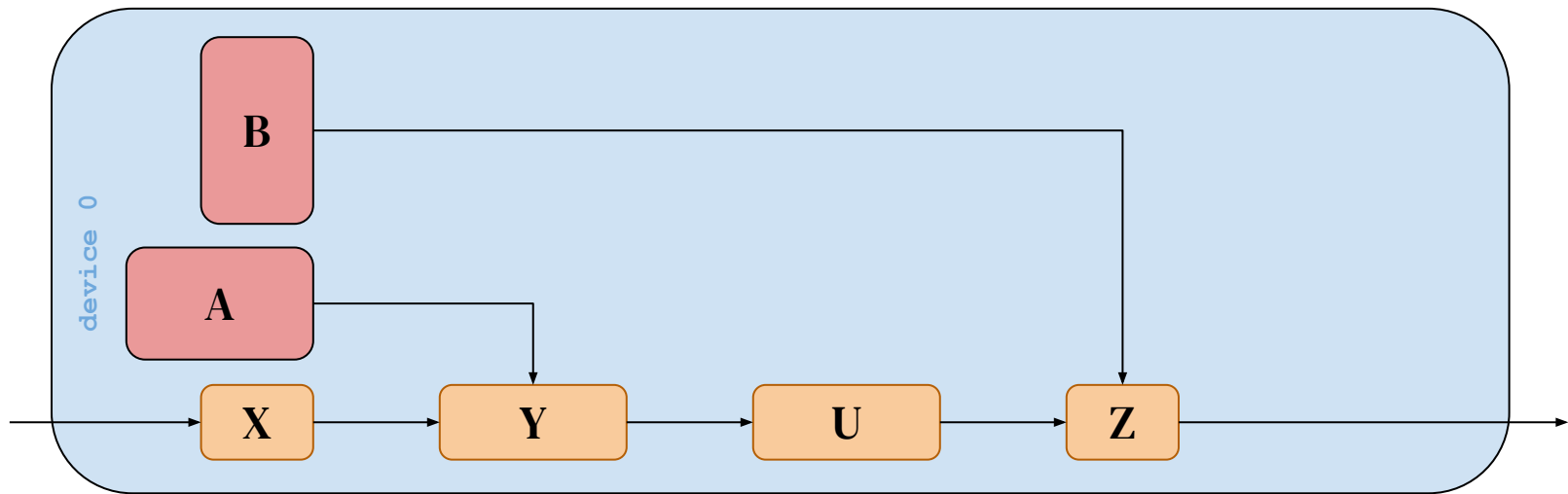
```
def mlp(  
    x: Float[Array, "B H_1"],  
    A: Float[Array, "H_1 H_2"],  
    B: Float[Array, "H_2 H_1"]) -> Float[Array, "B H_1"]:  
    """ Z = max(X * W_1, 0) * W_2 """  
    y = jnp.dot(x, A)      # [B, H_1] @ [H_1, H_2] -> [B, H_2]  
    u = jnp.maximum(y, 0)  # [B, H_2]  
    z = jnp.dot(u, B)      # [B, H_2] @ [H_2, H_1] -> [B, H_1]  
    return z
```



Sharding a 2-layer MLP

```
def mlp(  
    x: Float[Array, "B H_1"],  
    A: Float[Array, "H_1 H_2"],  
    B: Float[Array, "H_2 H_1"]) -> Float[Array, "B H_1"]:  
    """ Z = max(X * W_1, 0) * W_2 """  
    y = jnp.dot(x, A)      # [B, H_1] @ [H_1, H_2] -> [B, H_2]  
    u = jnp.maximum(y, 0)  # [B, H_2]  
    z = jnp.dot(u, B)      # [B, H_2] @ [H_2, H_1] -> [B, H_1]  
    return z
```

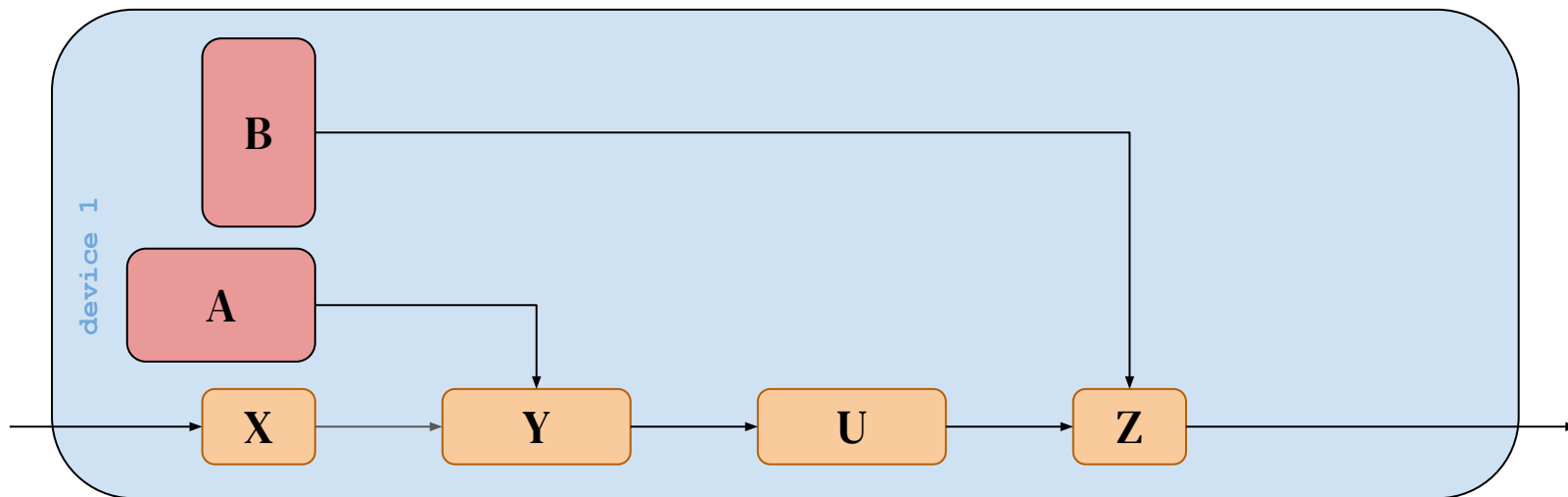
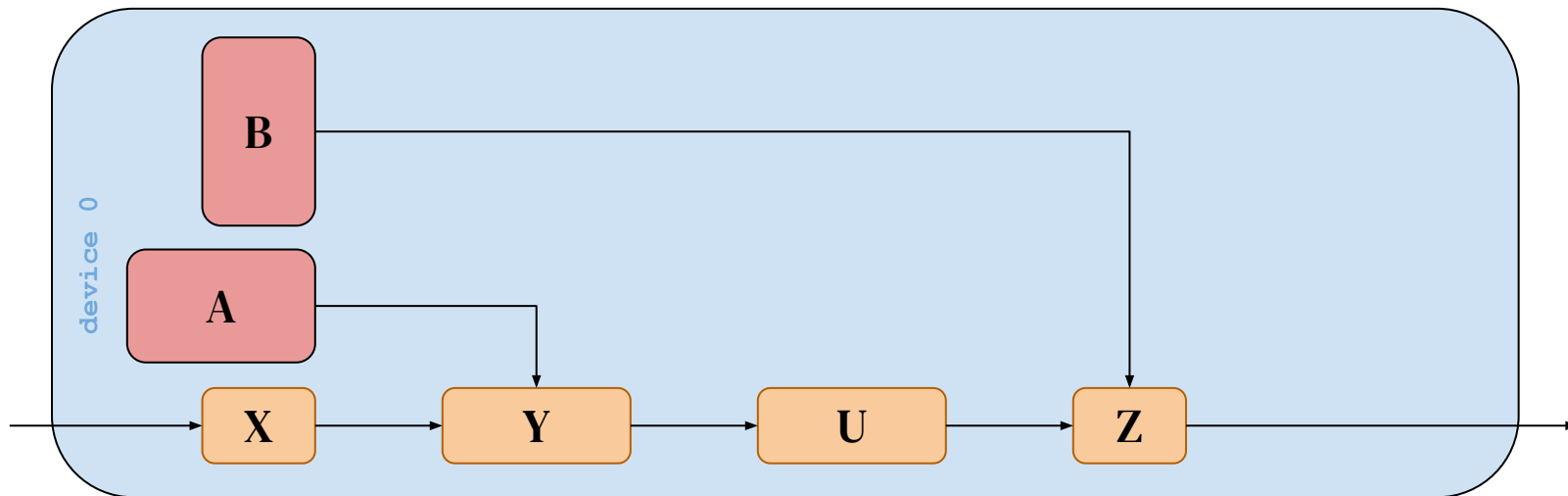




$B: [H_2,$
 $H_1]$

$A: [H_1,$
 $H_2]$

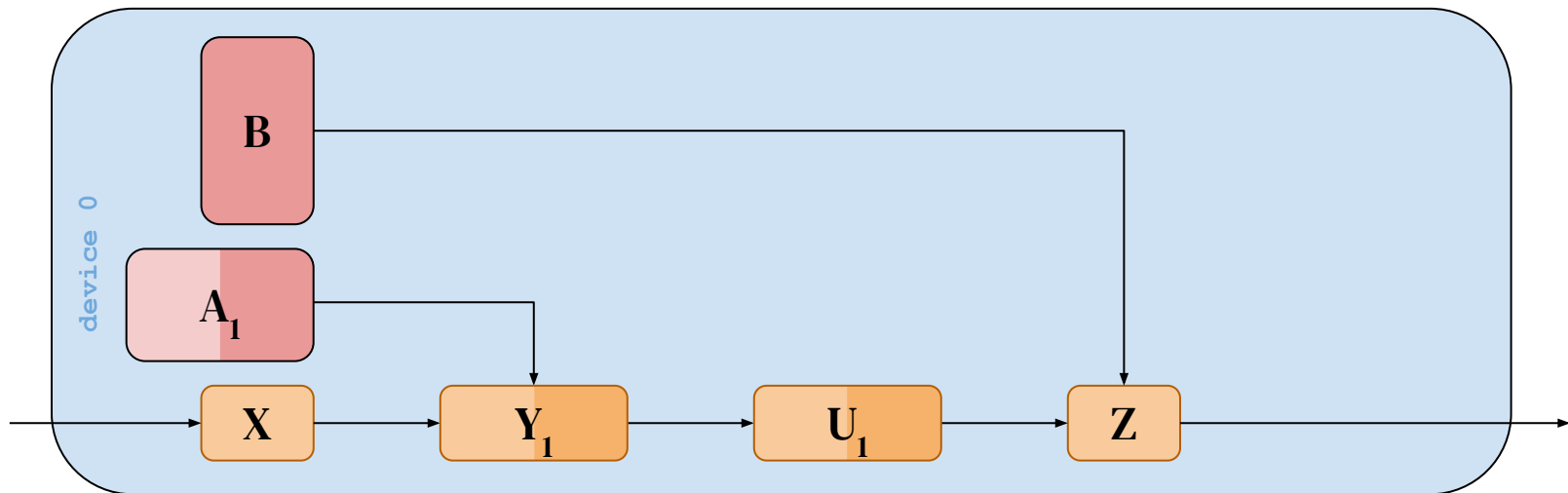
$X: [B, H_1]$



$B: [H_2, H_1]$

$A_1: [H_1,$
 $H_{2/2}]$

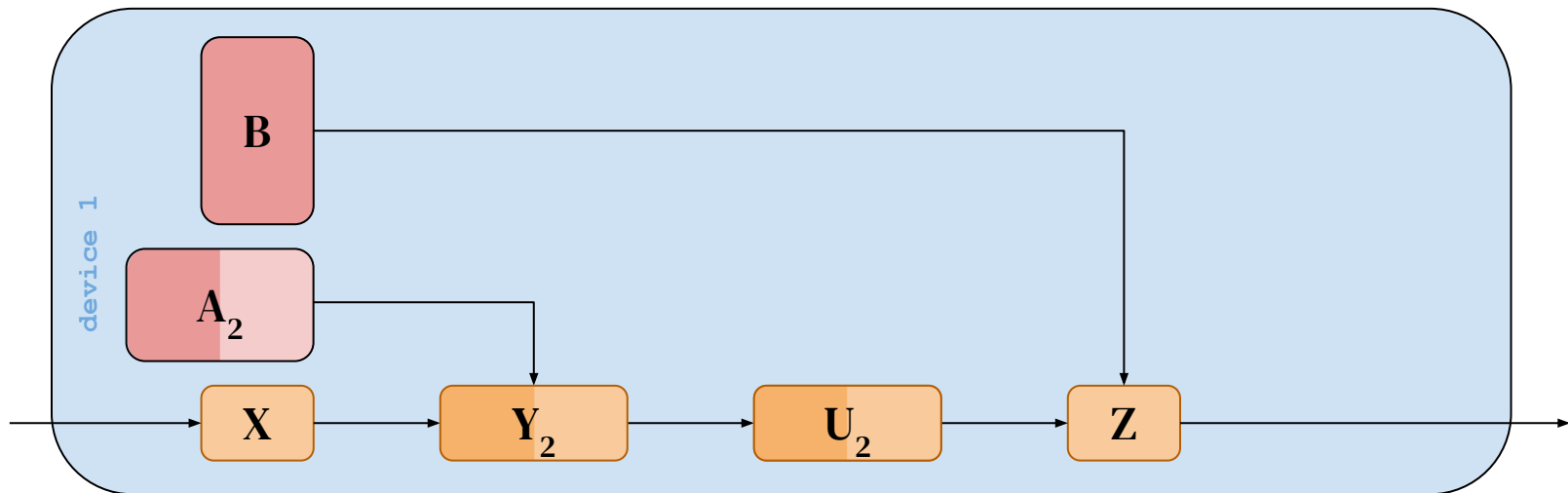
$X: [B, H_1]$



$B: [H_2, H_1]$

$A_2: [H_1,$
 $H_{2/2}]$

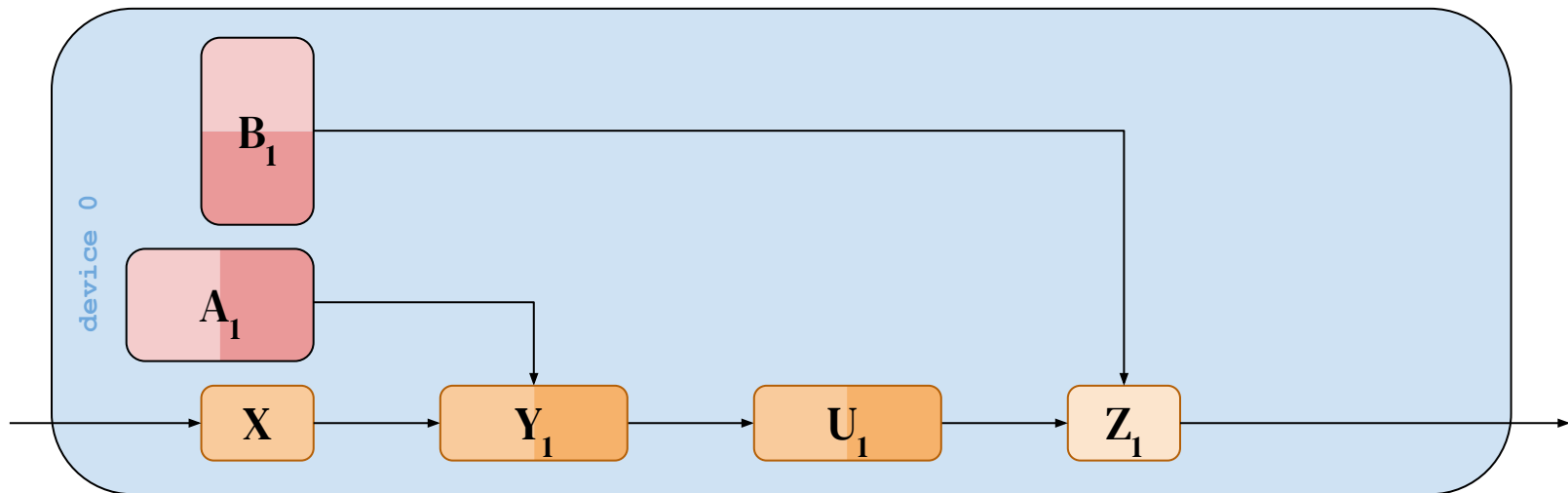
$X: [B, H_1]$



$B_1 : [H_{2/2},$
 $H_1]$

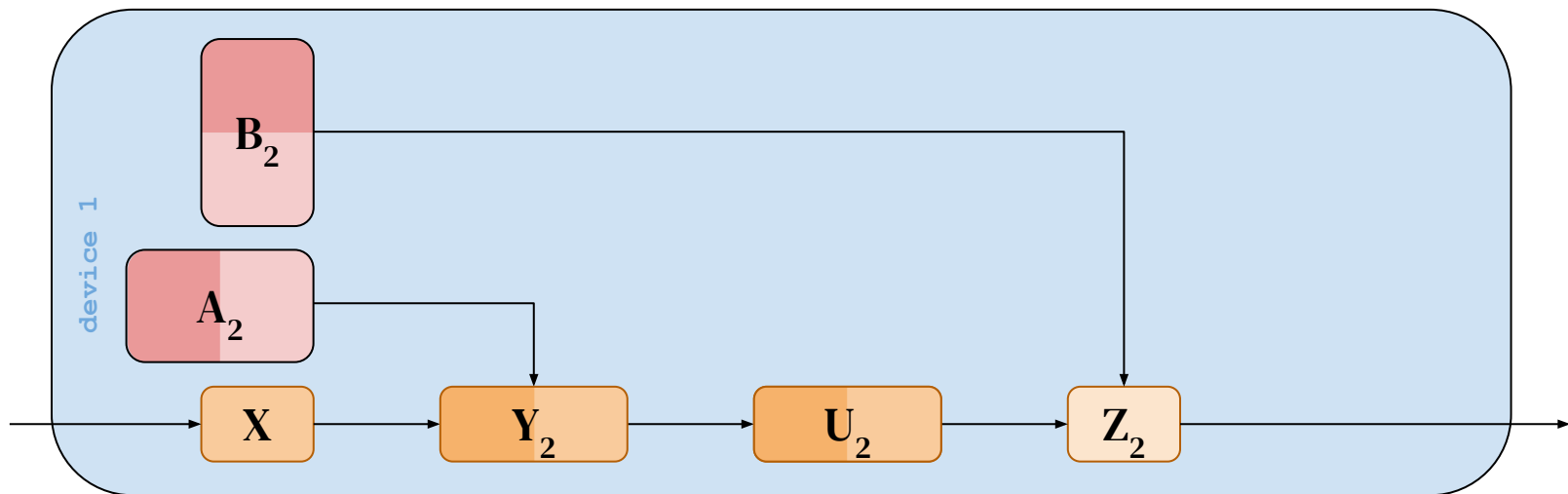
$A_1 : [H_1,$
 $H_{2/2}]$

$X : [B, H_1]$



$B_2 : [H_{2/2},$
 $H_1]$

$A_2 : [H_1,$
 $H_{2/2}]$



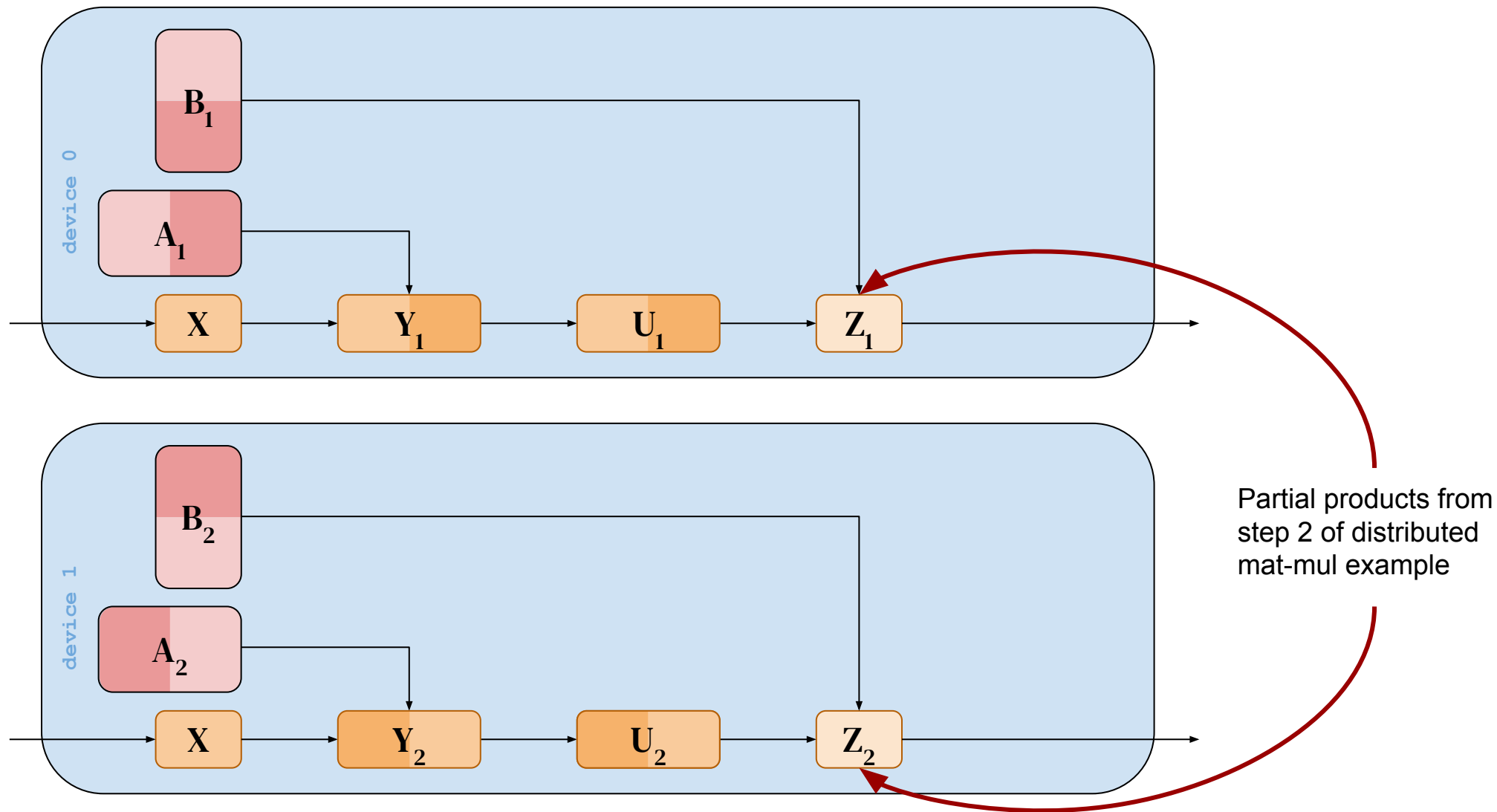
$\mathbf{B}_1 : [\mathbf{H}_{2/2},$
 $\mathbf{H}_1]$

$\mathbf{A}_1 : [\mathbf{H}_1,$
 $\mathbf{H}_{2/2}]$

$\mathbf{X} : [\mathbf{B}, \mathbf{H}_1]$

$\mathbf{B}_2 : [\mathbf{H}_{2/2},$
 $\mathbf{H}_1]$

$\mathbf{A}_2 : [\mathbf{H}_1,$
 $\mathbf{H}_{2/2}]$



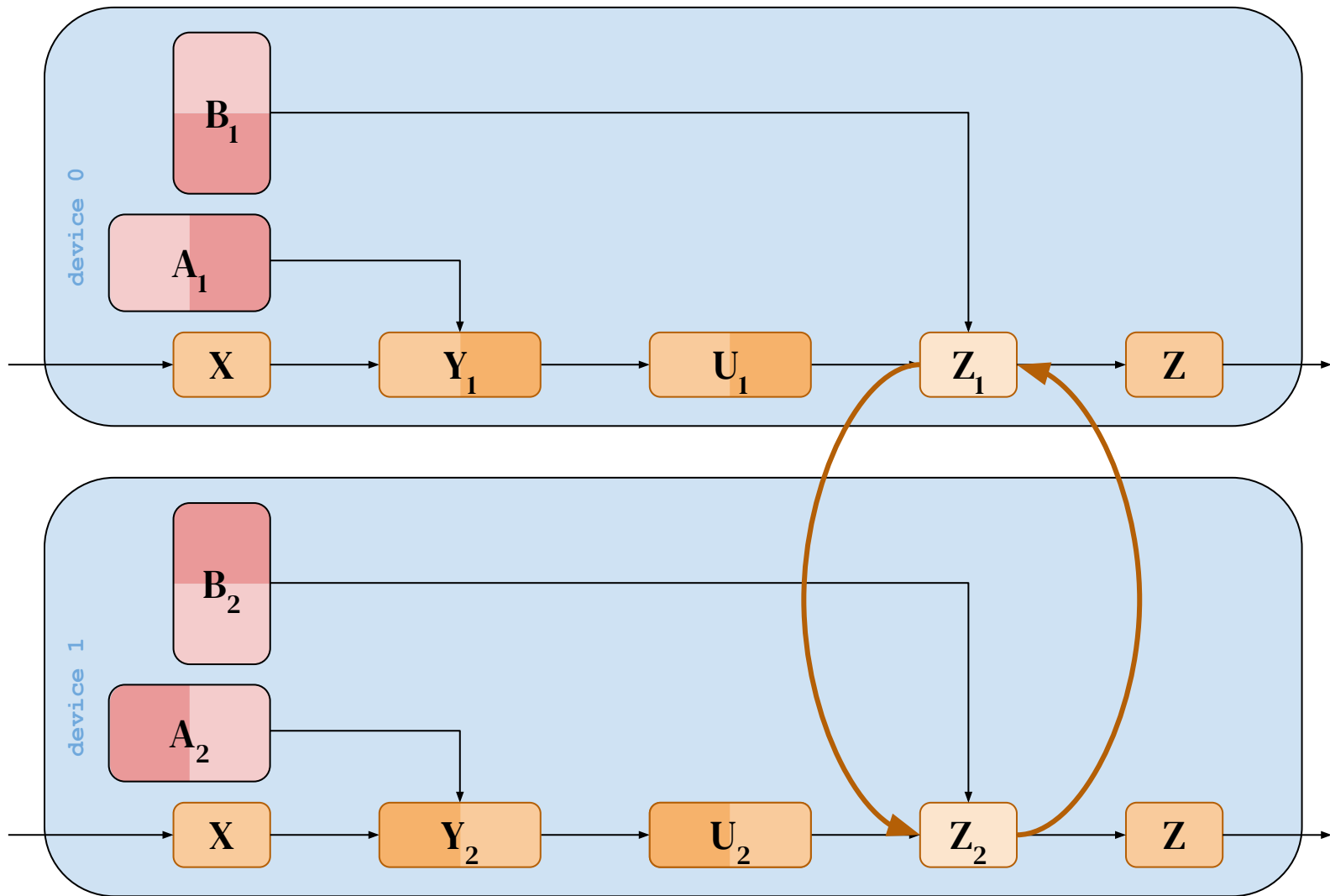
$B_1 : [H_{2/2},$
 $H_1]$

$A_1 : [H_1,$
 $H_{2/2}]$

$X : [B, H_1]$

$B_2 : [H_{2/2},$
 $H_1]$

$A_2 : [H_1,$
 $H_{2/2}]$



AllReduce

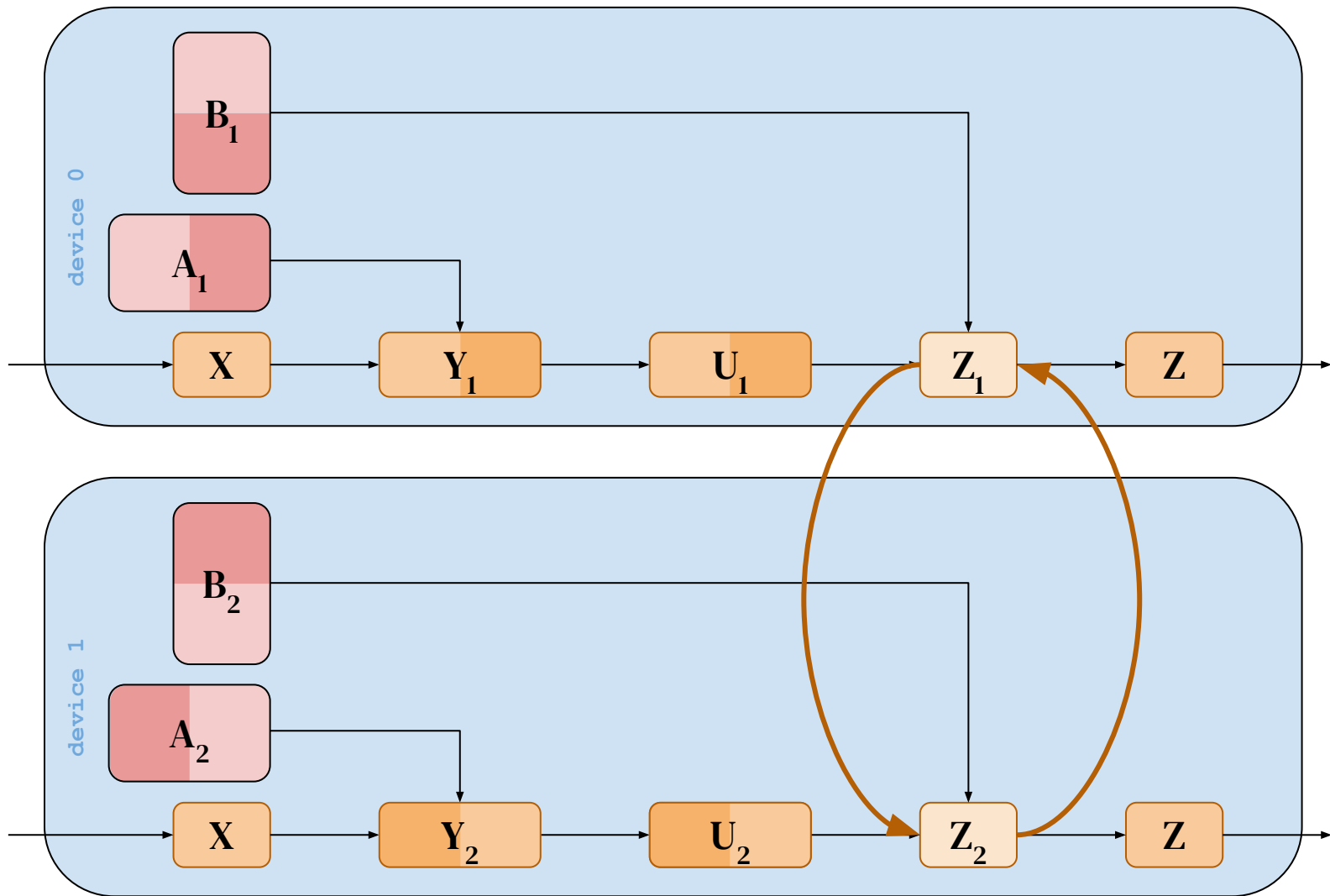
$B_1 : [H_{2/2}, H_1]$

$A_1 : [H_1, H_{2/2}]$

$X : [B, H_1]$

$B_2 : [H_{2/2}, H_1]$

$A_2 : [H_1, H_{2/2}]$



What about gradients?

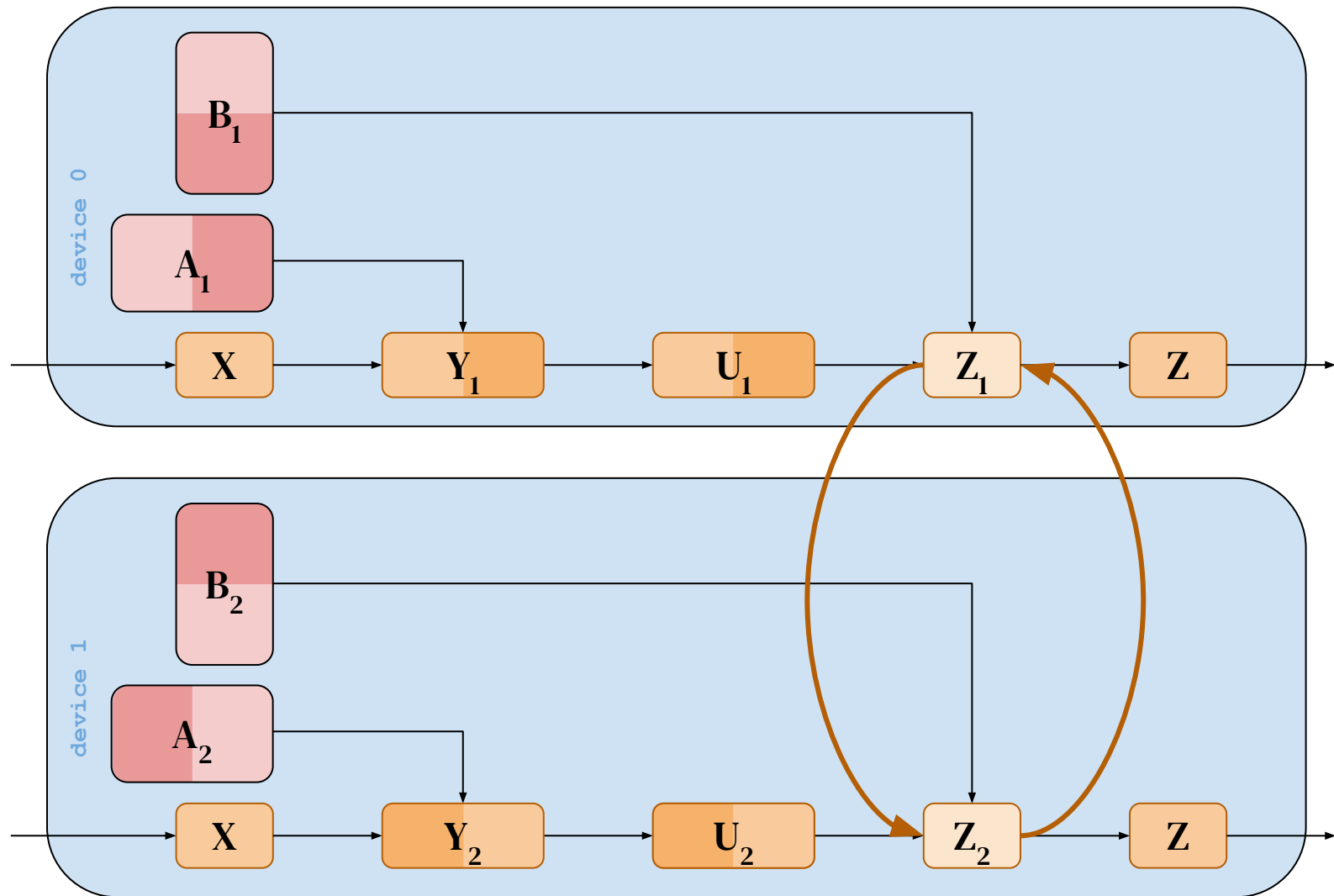
$B_1 : [H_{2/2}, H_1]$

$A_1 : [H_1, H_{2/2}]$

$X : [B, H_1]$

$B_2 : [H_{2/2}, H_1]$

$A_2 : [H_1, H_{2/2}]$



Let's forget
about
devices for a
while

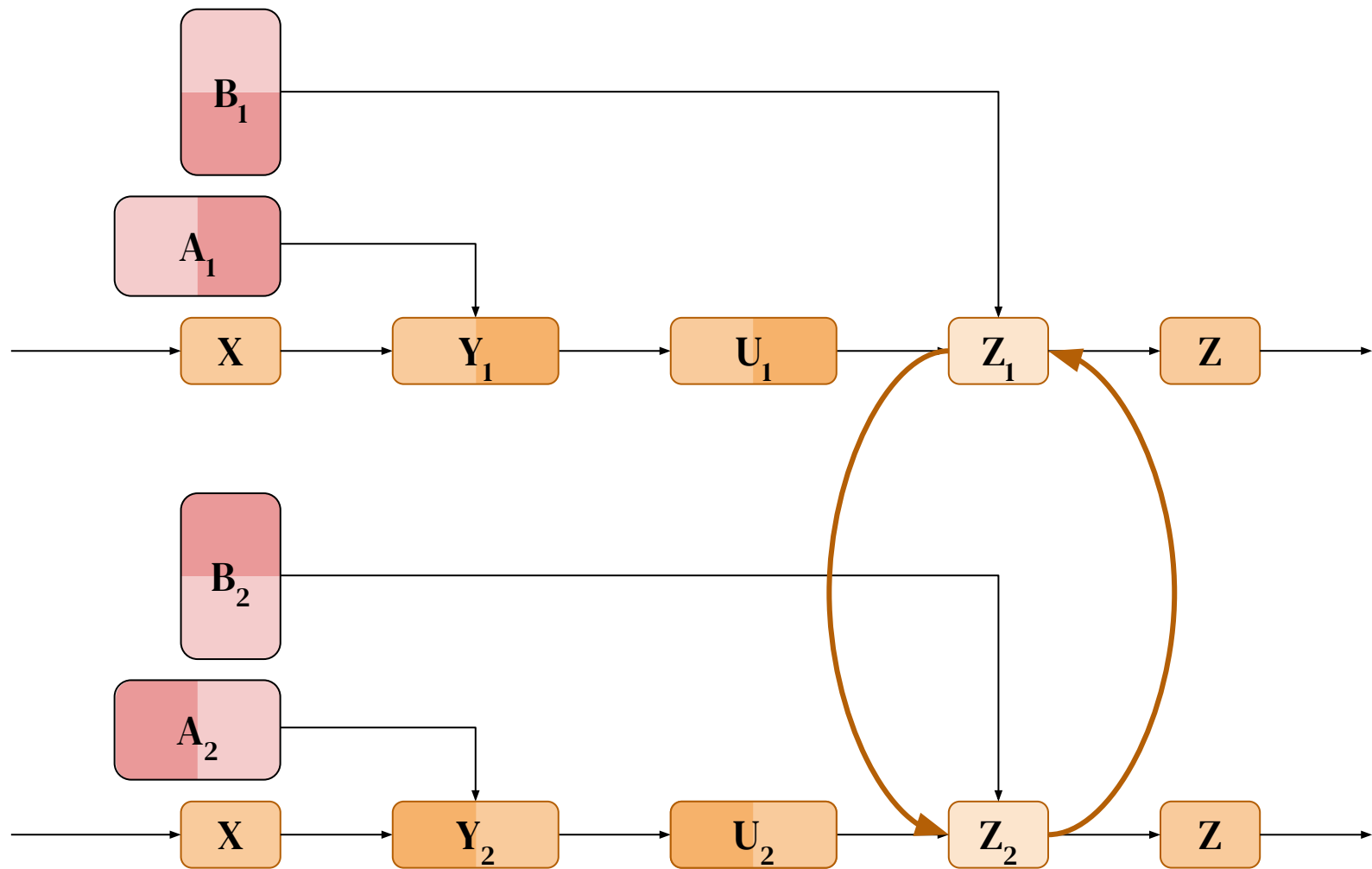
$B_1 : [H_{2/2}, H_1]$

$A_1 : [H_1, H_{2/2}]$

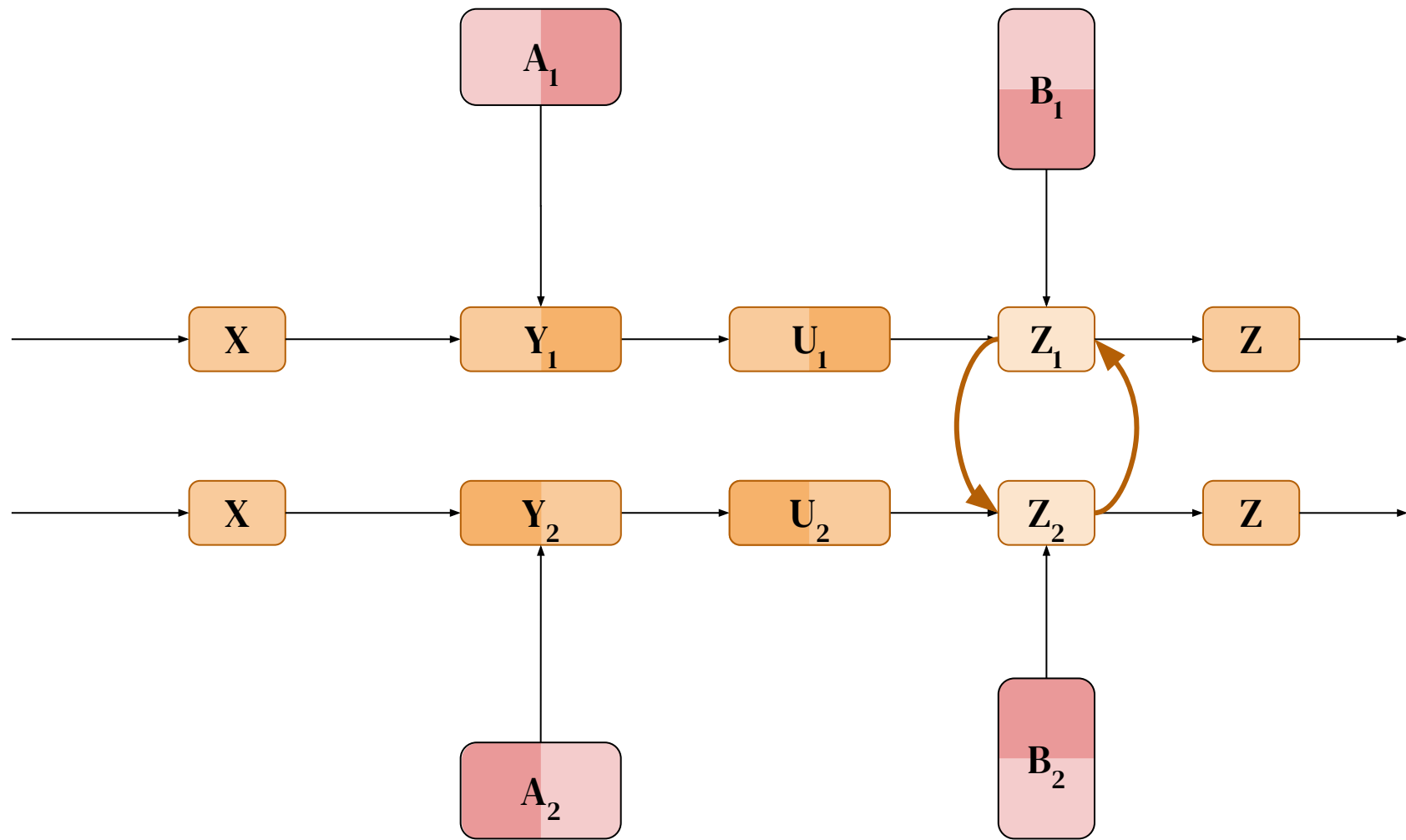
$X : [B, H_1]$

$B_2 : [H_{2/2}, H_1]$

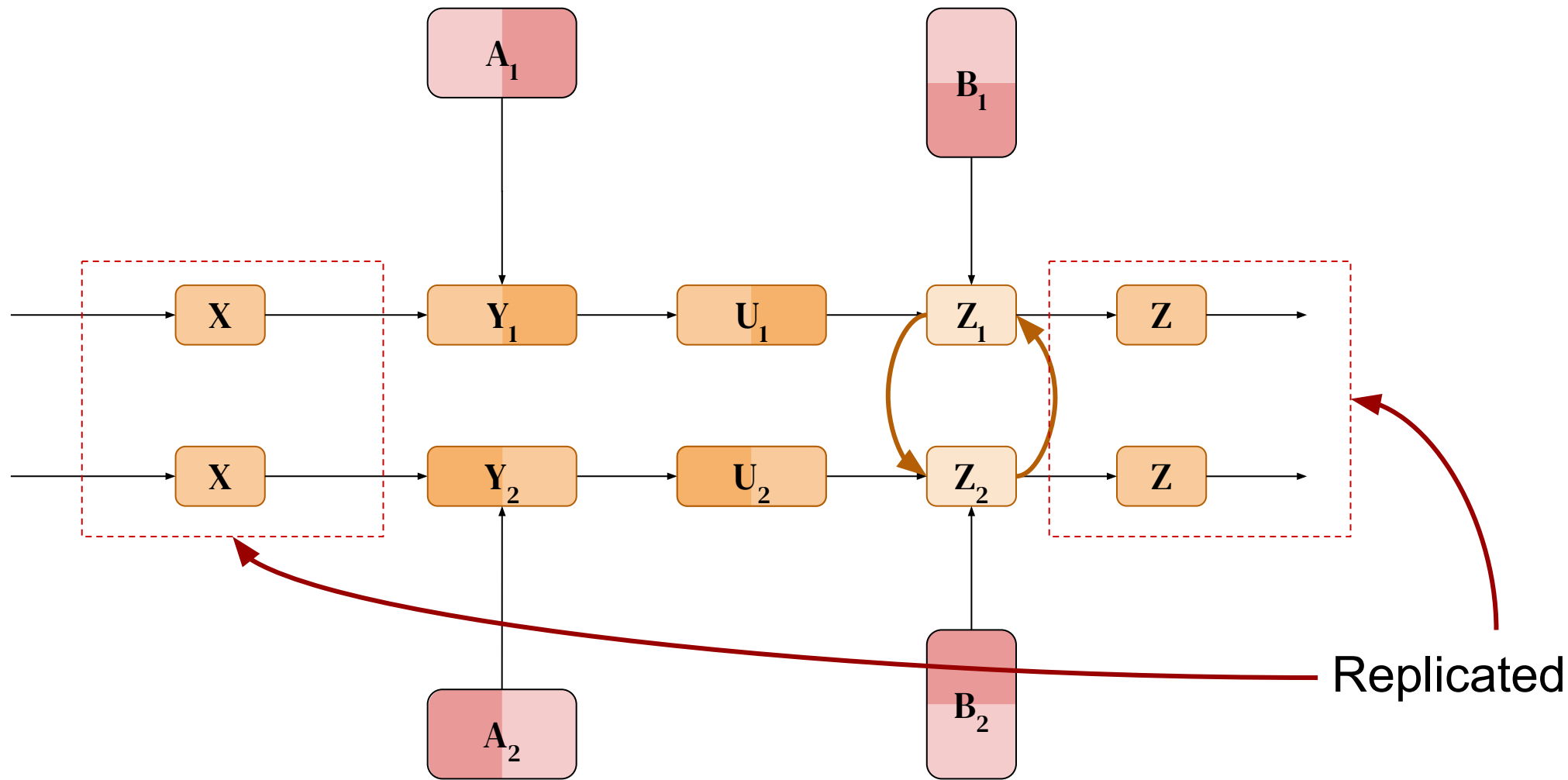
$A_2 : [H_1, H_{2/2}]$

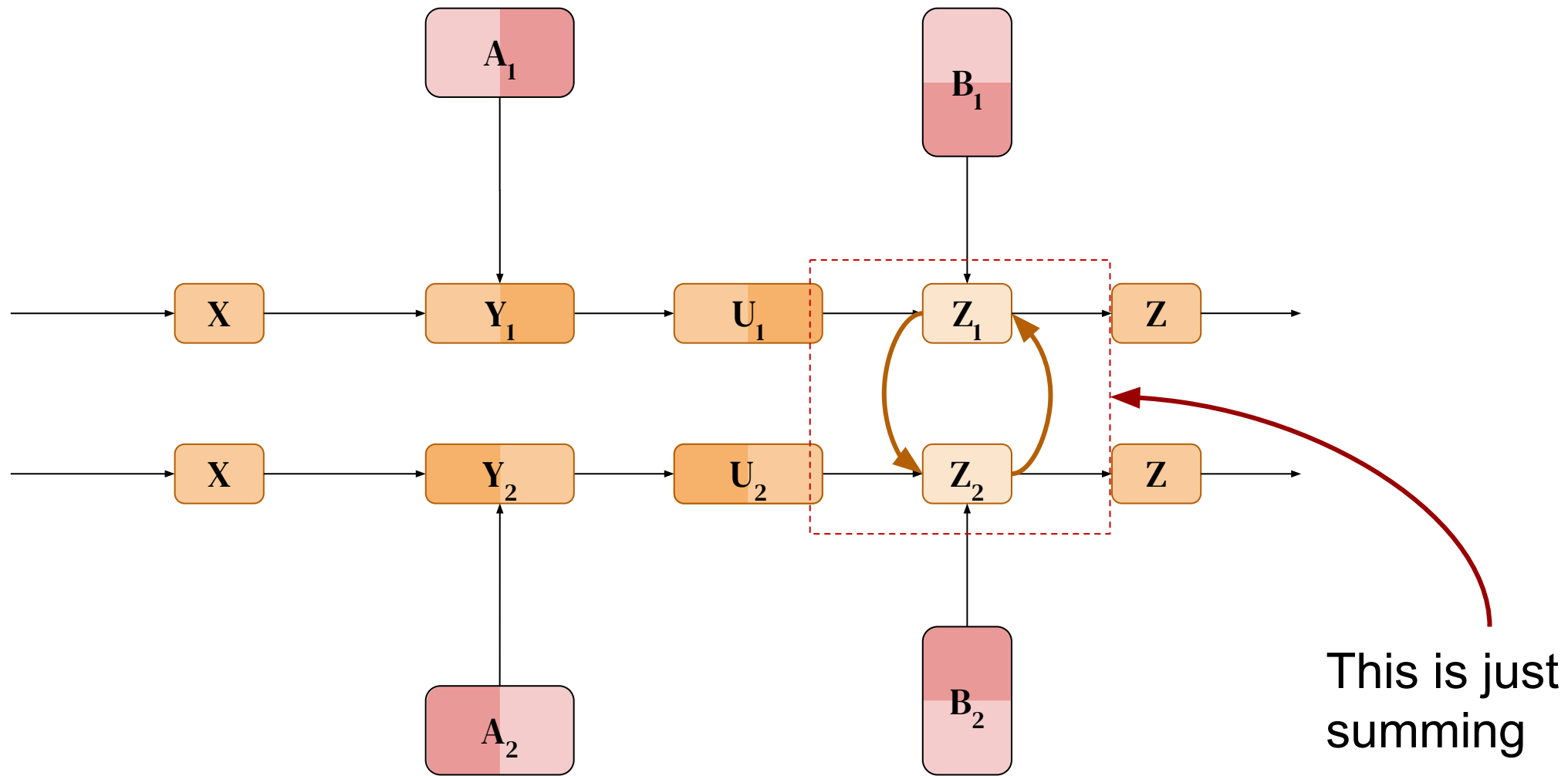


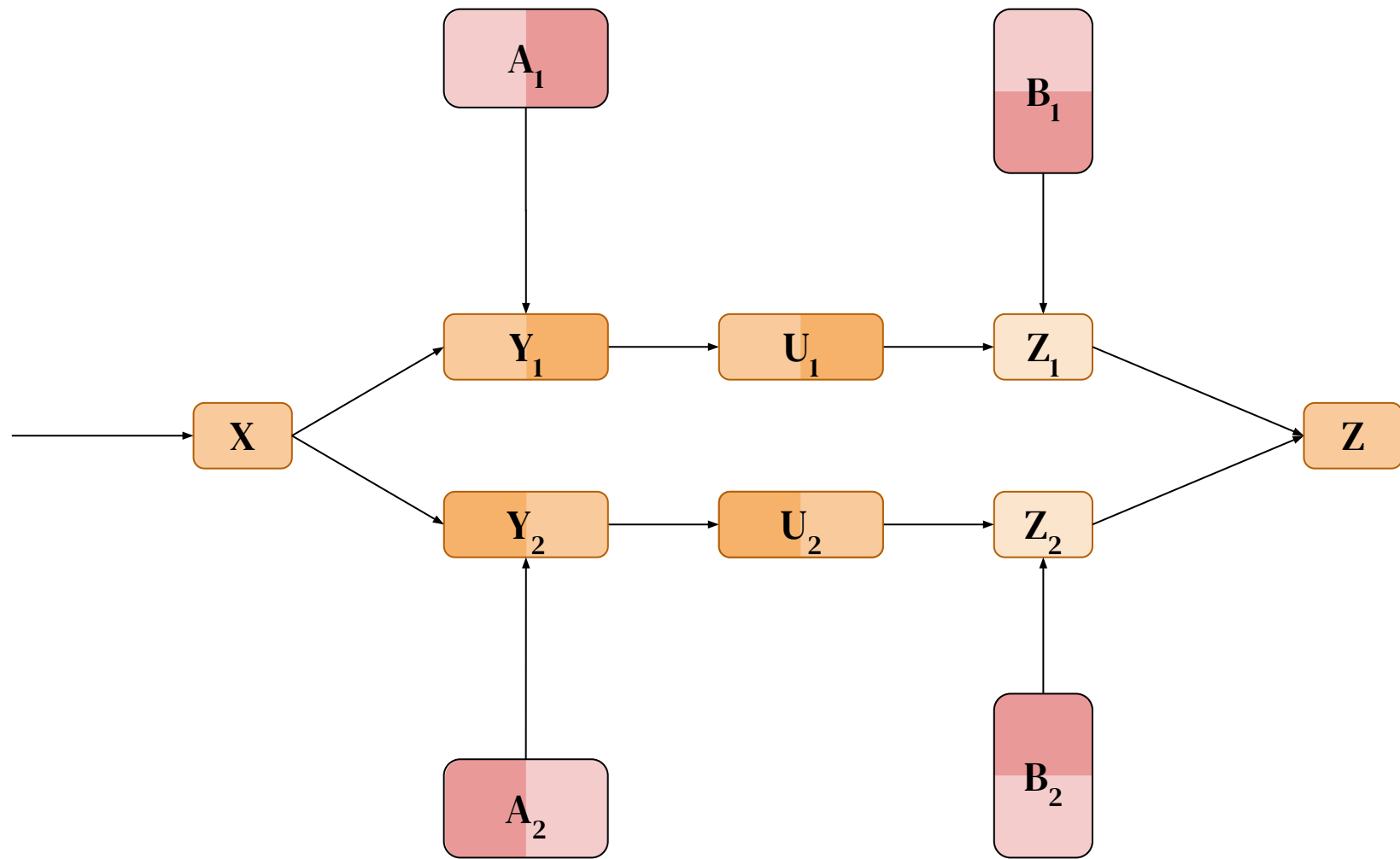
Let's forget
about
devices for a
while



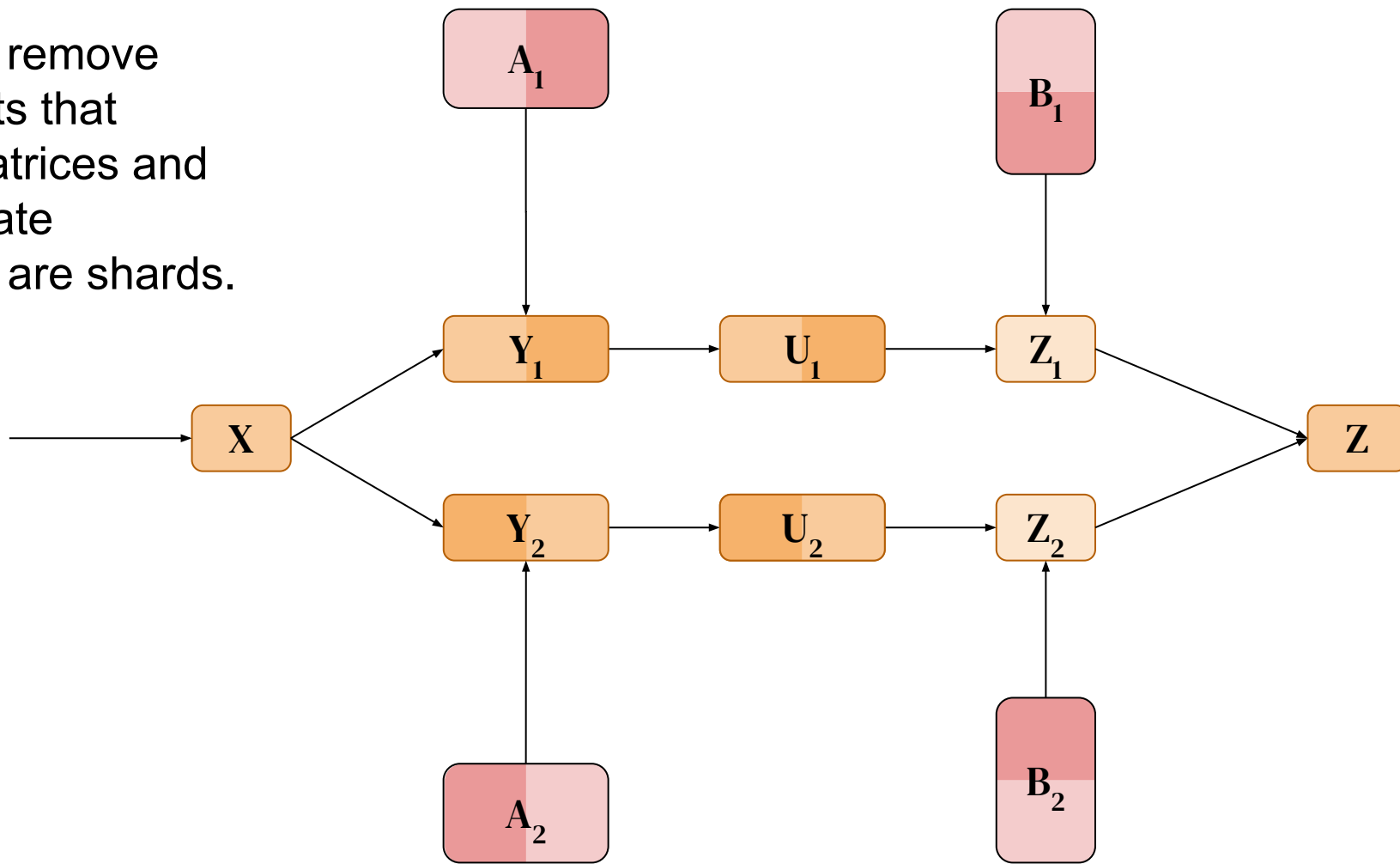
And rearrange
the graph a bit





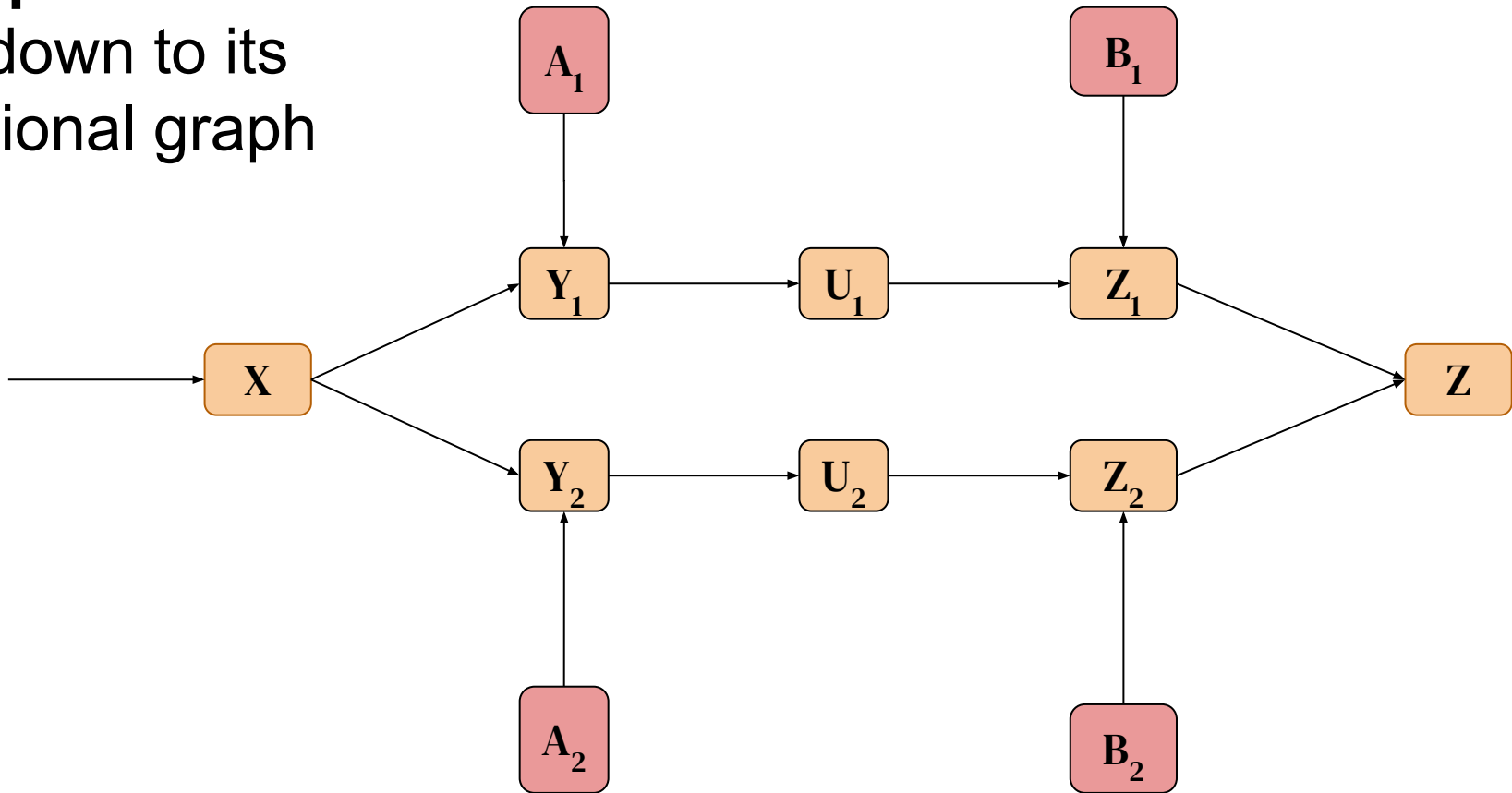


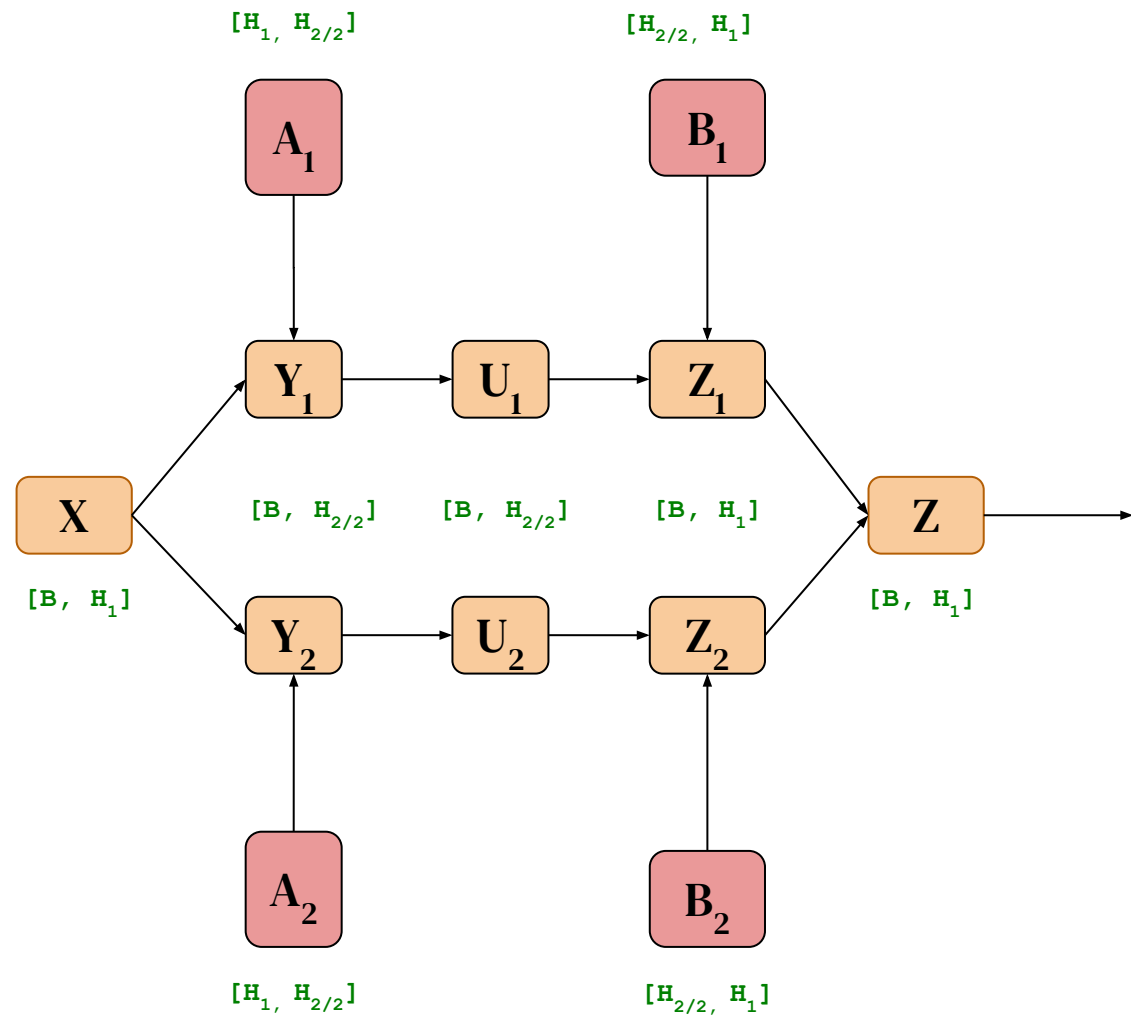
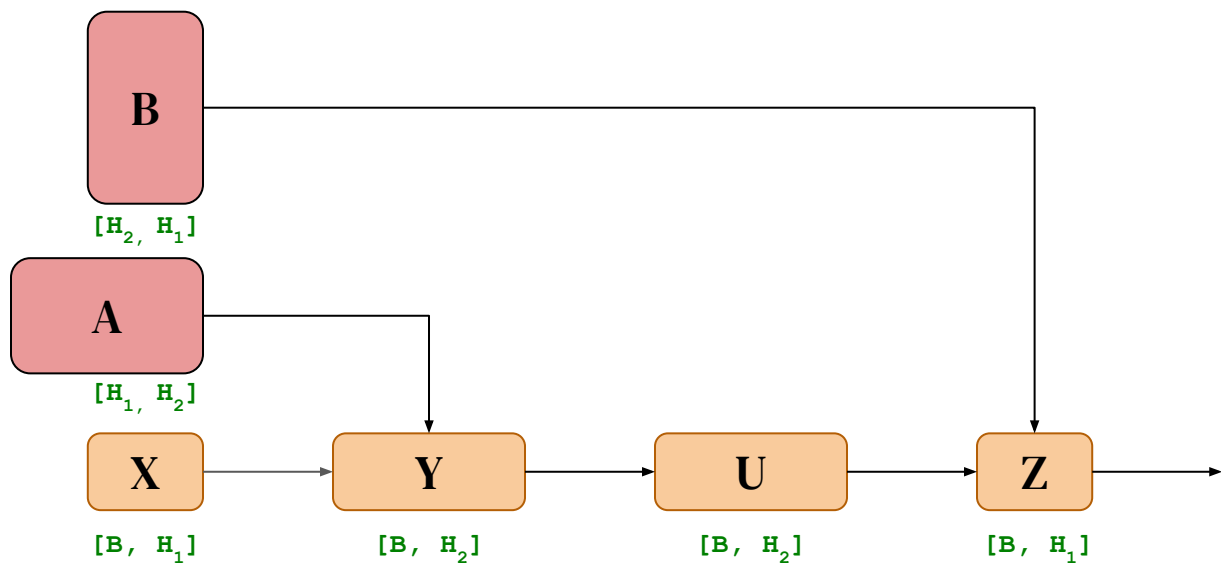
Let's also remove
visual hints that
weight matrices and
intermediate
activation are shards.



Forward pass

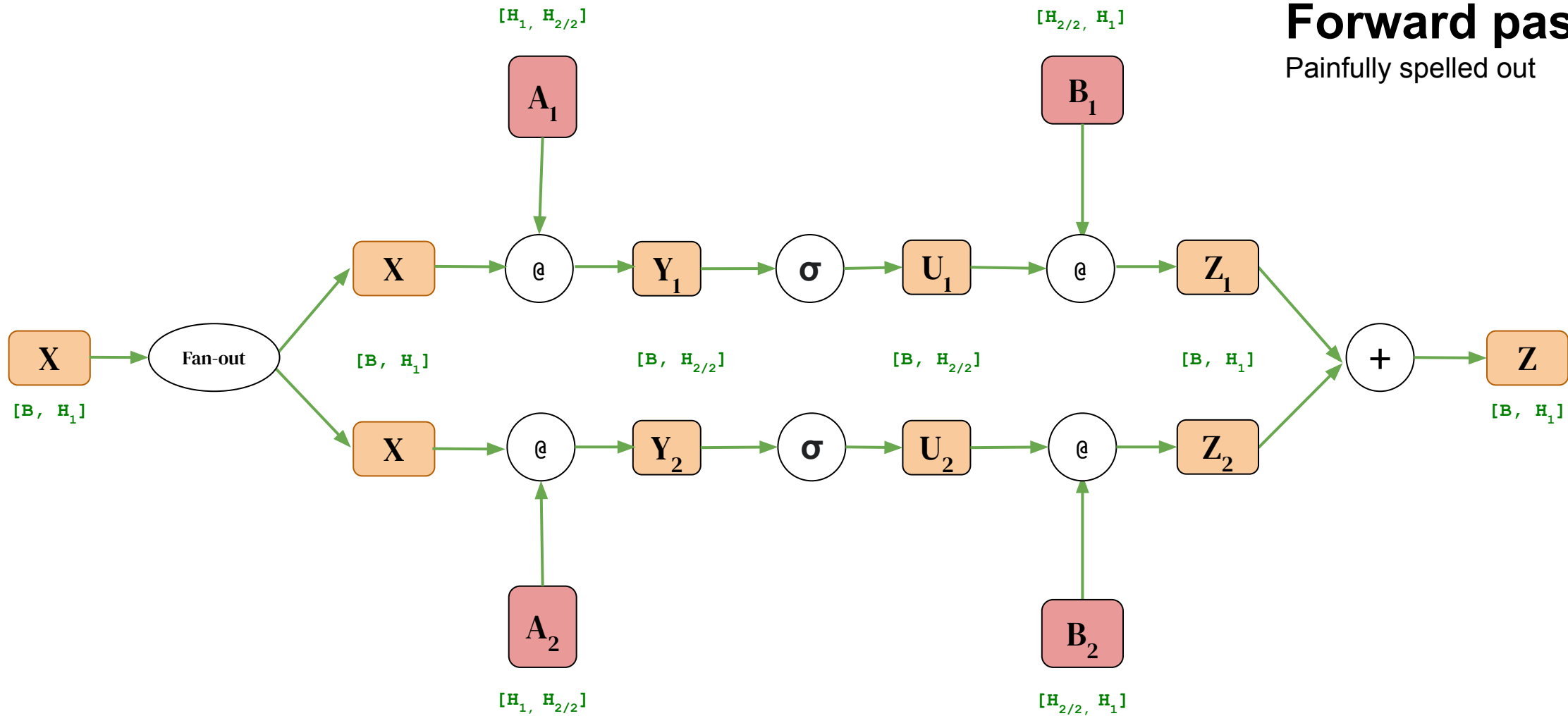
stripped down to its
computational graph





Forward pass

Two views on the computations in the forward pass.

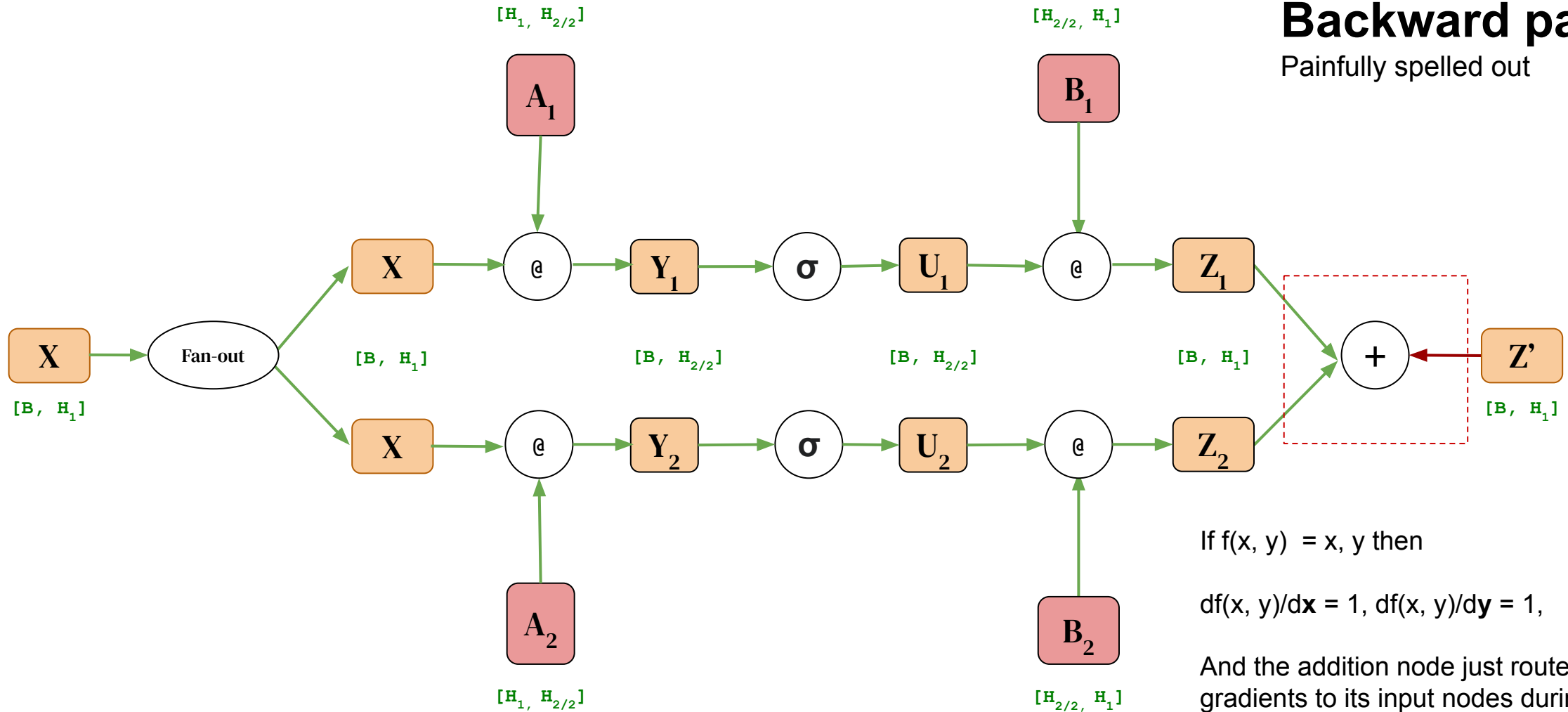


Forward pass

Painfully spelled out

Backward pass

Painfully spelled out



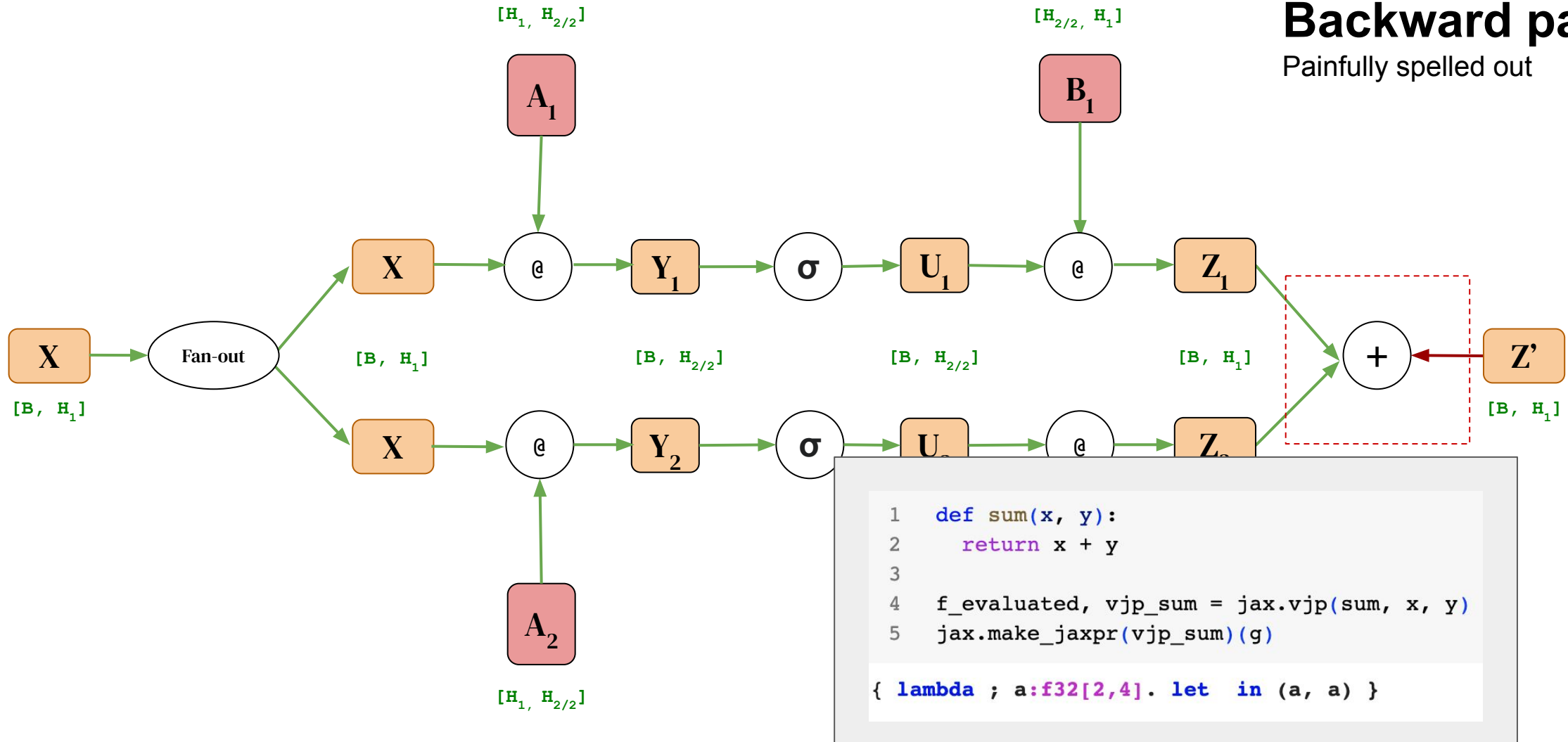
If $f(x, y) = x, y$ then

$df(x, y)/dx = 1, df(x, y)/dy = 1,$

And the addition node just routes gradients to its input nodes during backward pass

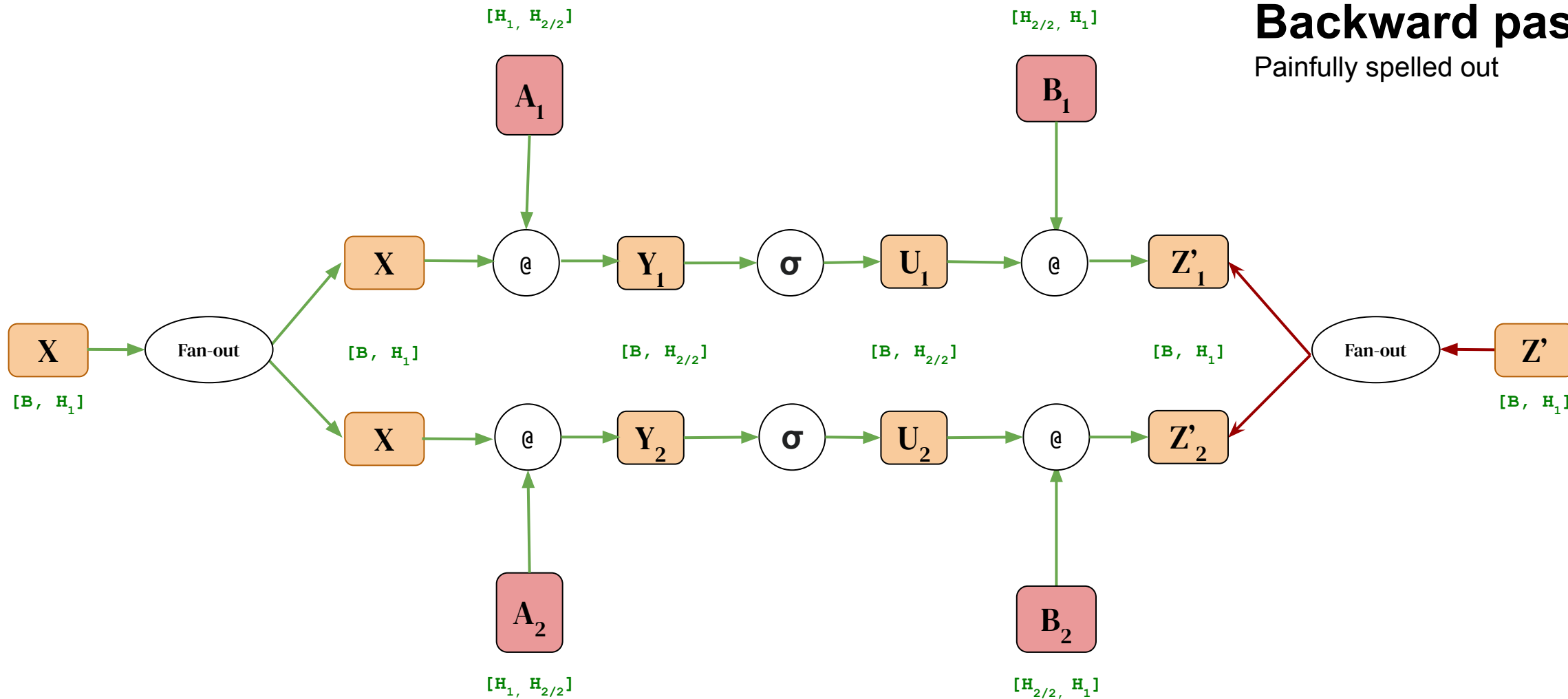
Backward pass

Painfully spelled out



Backward pass

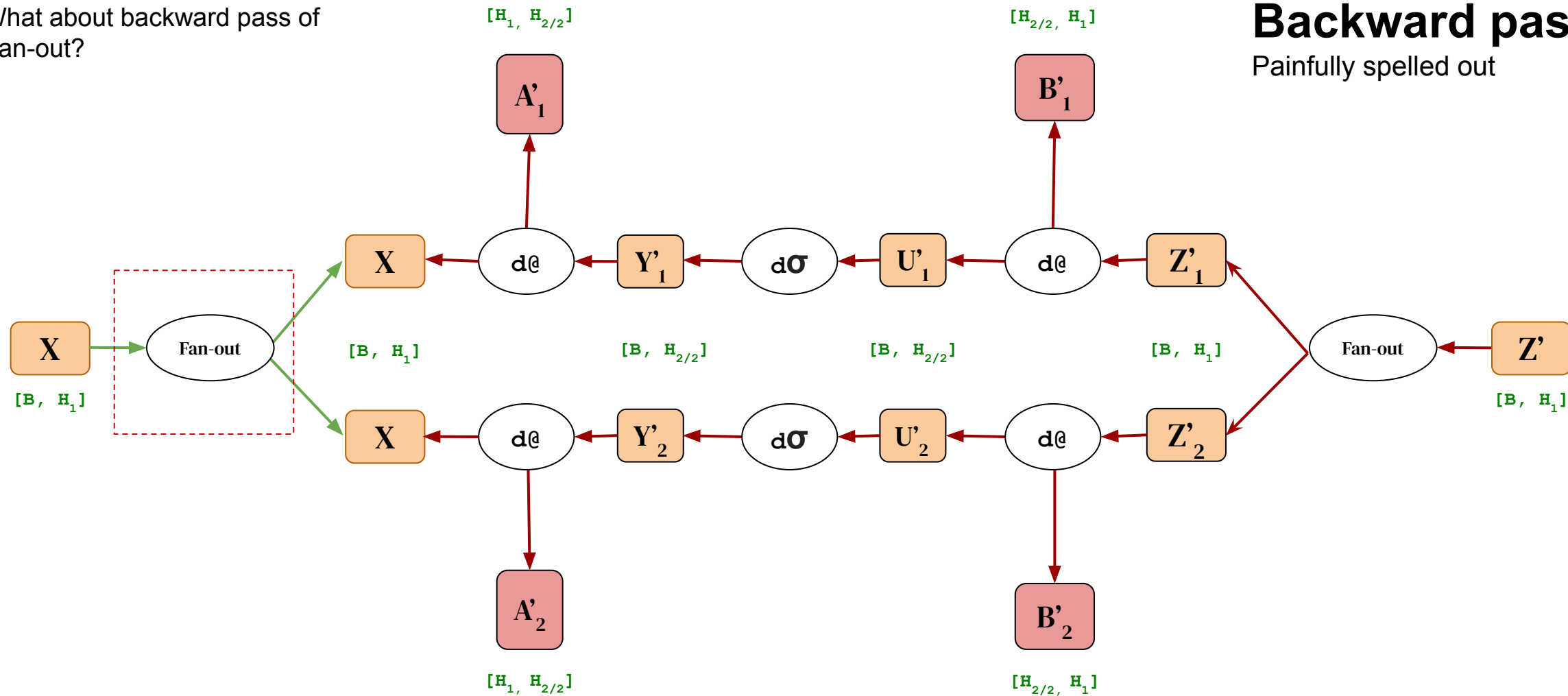
Painfully spelled out



What about backward pass of Fan-out?

Backward pass

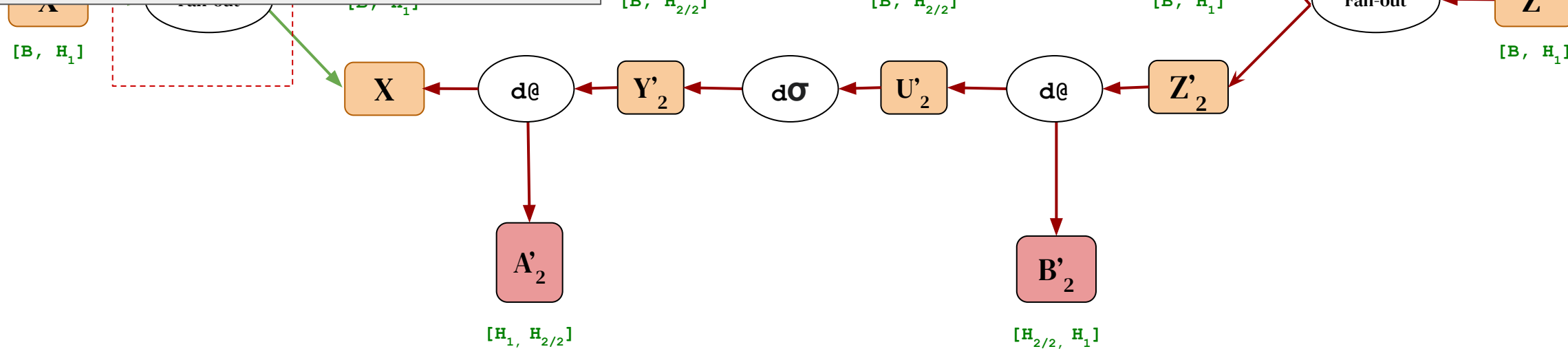
Painfully spelled out



$$\text{Fan-out}(X) = \begin{bmatrix} X \\ X \end{bmatrix} = \begin{bmatrix} I \\ I \end{bmatrix} X$$

Backward pass

Painfully spelled out

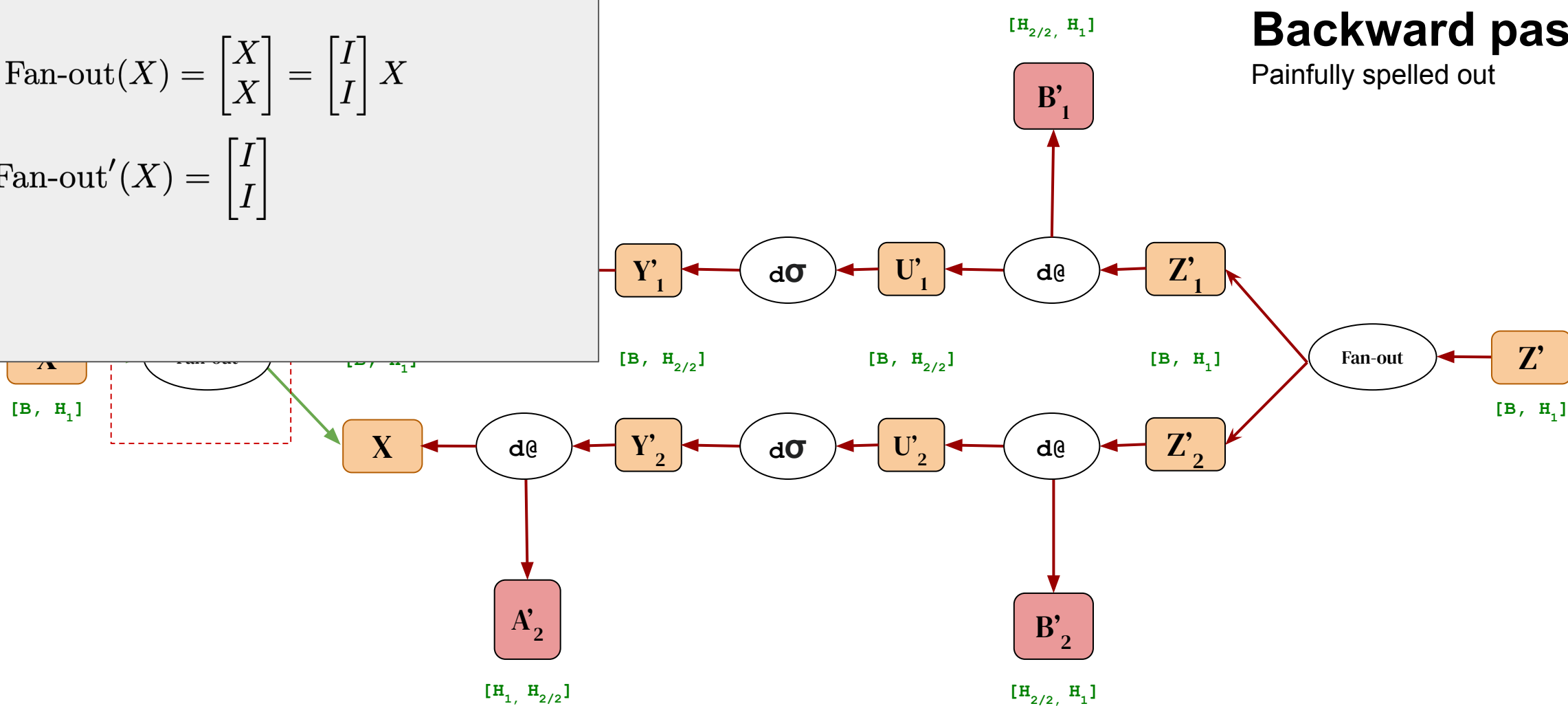


$$\text{Fan-out}(X) = \begin{bmatrix} X \\ X \end{bmatrix} = \begin{bmatrix} I \\ I \end{bmatrix} X$$

$$\text{Fan-out}'(X) = \begin{bmatrix} I \\ I \end{bmatrix} X$$

Backward pass

Painfully spelled out



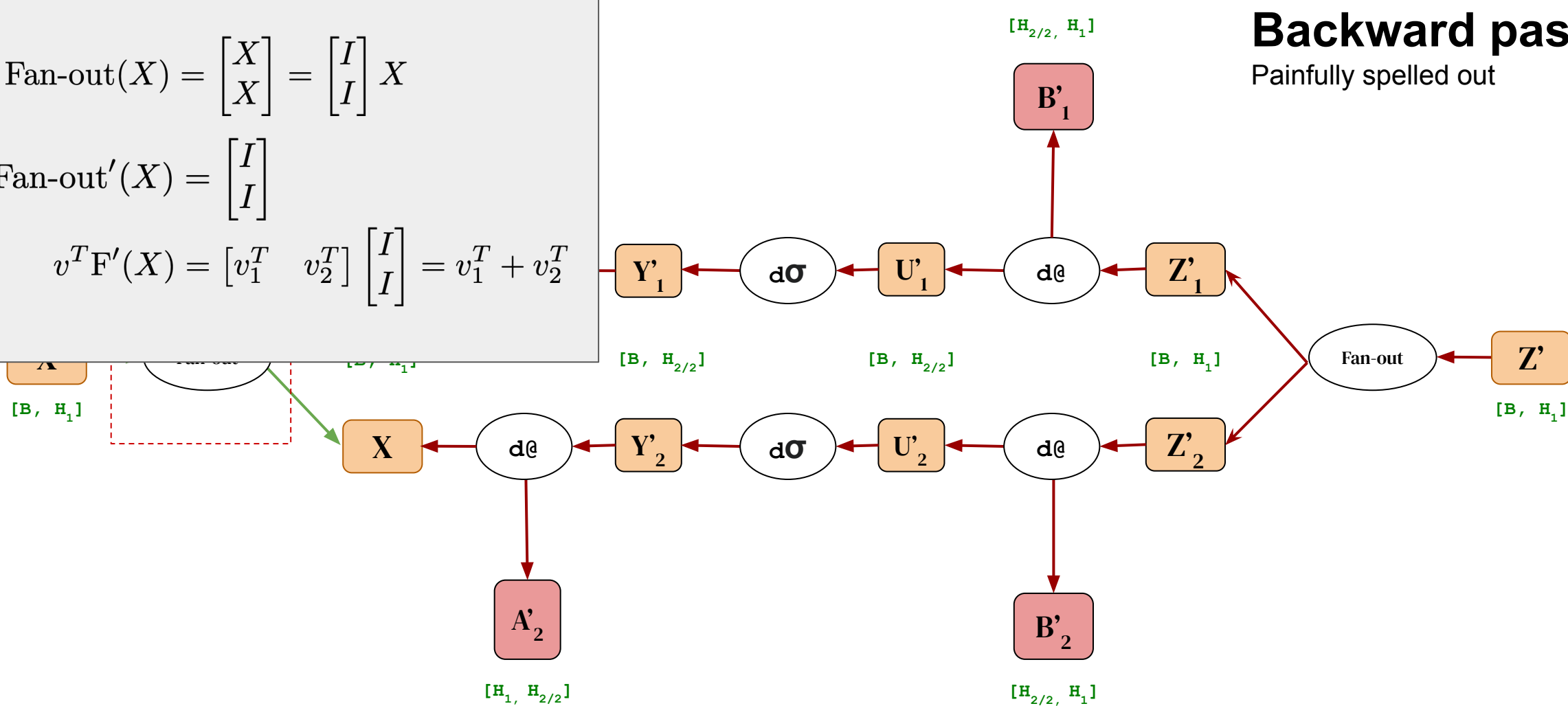
$$\text{Fan-out}(X) = \begin{bmatrix} X \\ X \end{bmatrix} = \begin{bmatrix} I \\ I \end{bmatrix} X$$

$$\text{Fan-out}'(X) = \begin{bmatrix} I \\ I \end{bmatrix}$$

$$v^T F'(X) = \begin{bmatrix} v_1^T & v_2^T \end{bmatrix} \begin{bmatrix} I \\ I \end{bmatrix} = v_1^T + v_2^T$$

Backward pass

Painfully spelled out



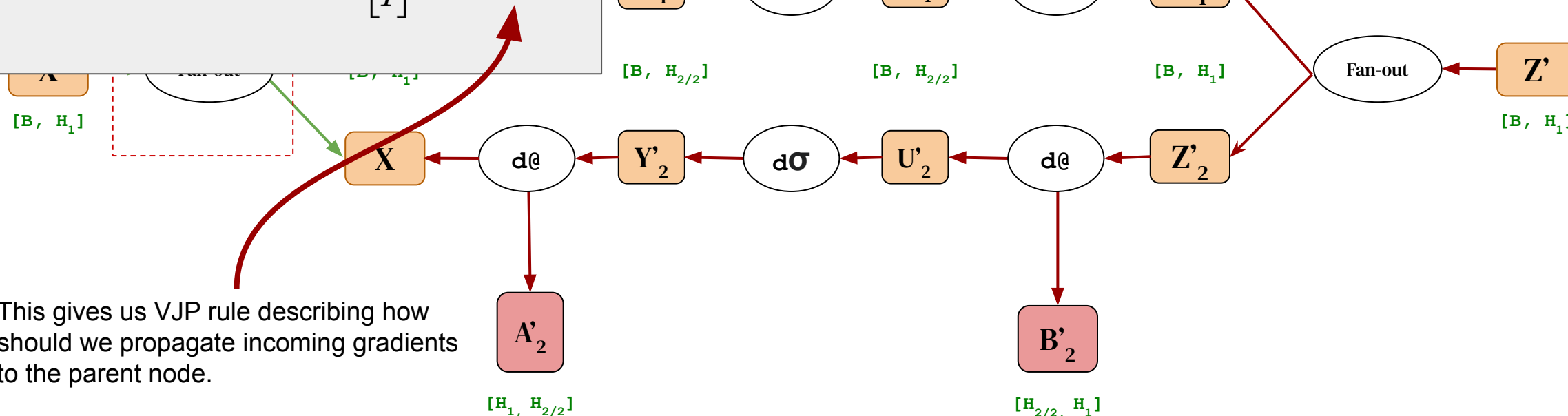
$$\text{Fan-out}(X) = \begin{bmatrix} X \\ X \end{bmatrix} = \begin{bmatrix} I \\ I \end{bmatrix} X$$

$$\text{Fan-out}'(X) = \begin{bmatrix} I \\ I \end{bmatrix}$$

$$v^T F'(X) = \begin{bmatrix} v_1^T & v_2^T \end{bmatrix} \begin{bmatrix} I \\ I \end{bmatrix} = v_1^T + v_2^T$$

Backward pass

Painfully spelled out



This gives us VJP rule describing how should we propagate incoming gradients to the parent node.

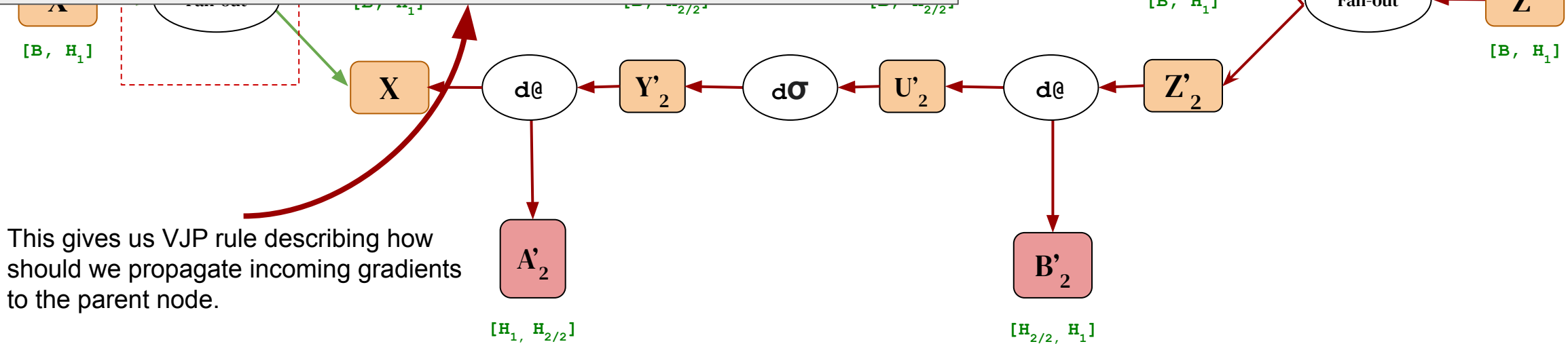
Backward pass

Painfully spelled out

```
1 def fan_out(x):
2     return [x, x]
3
4 f_evaluated, vjp_fan_out = jax.vjp(fan_out, x)
5 jax.make_jaxpr(vjp_fan_out)([g, g])
```



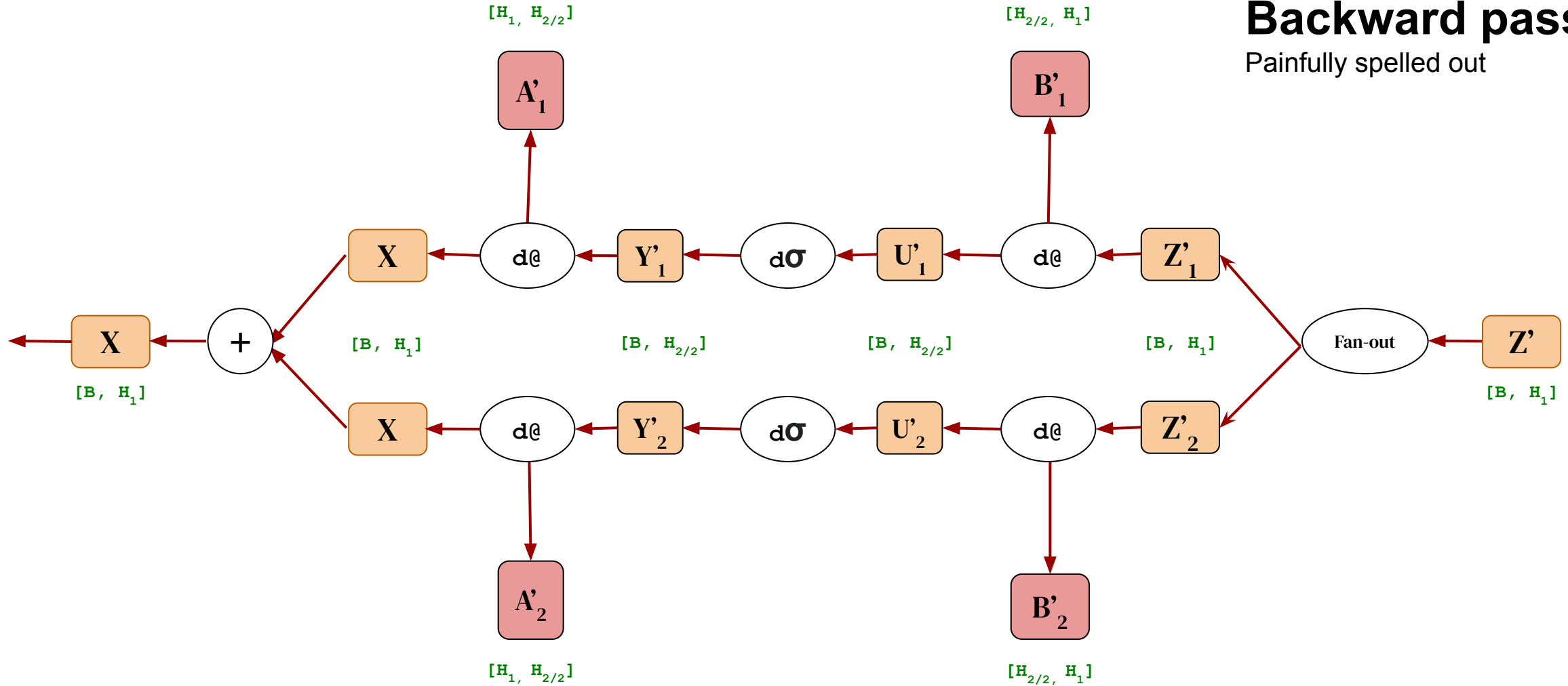
```
{ lambda ; a:f32[2,4] b:f32[2,4]. let c:f32[2,4] = add_any a b in (c,) }
```

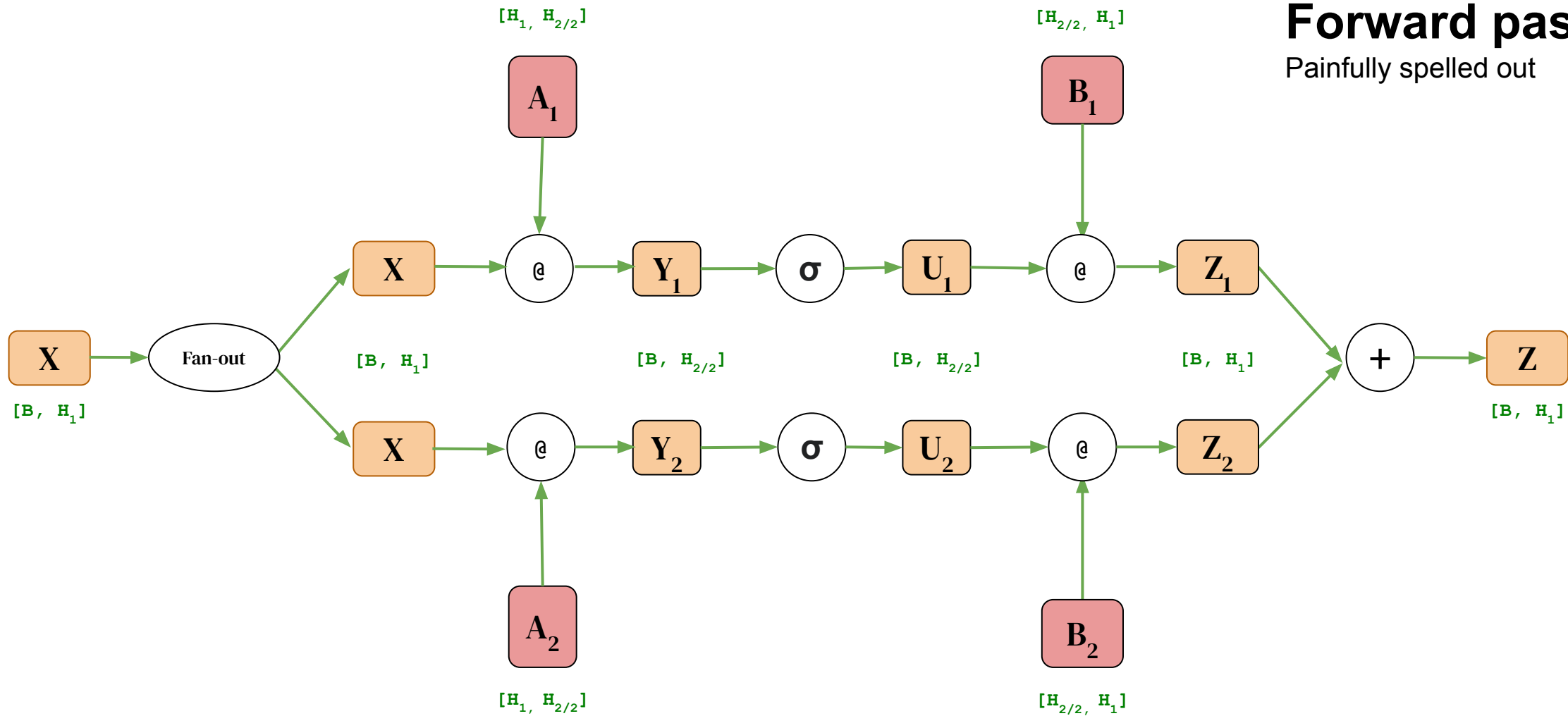


This gives us VJP rule describing how should we propagate incoming gradients to the parent node.

Backward pass

Painfully spelled out

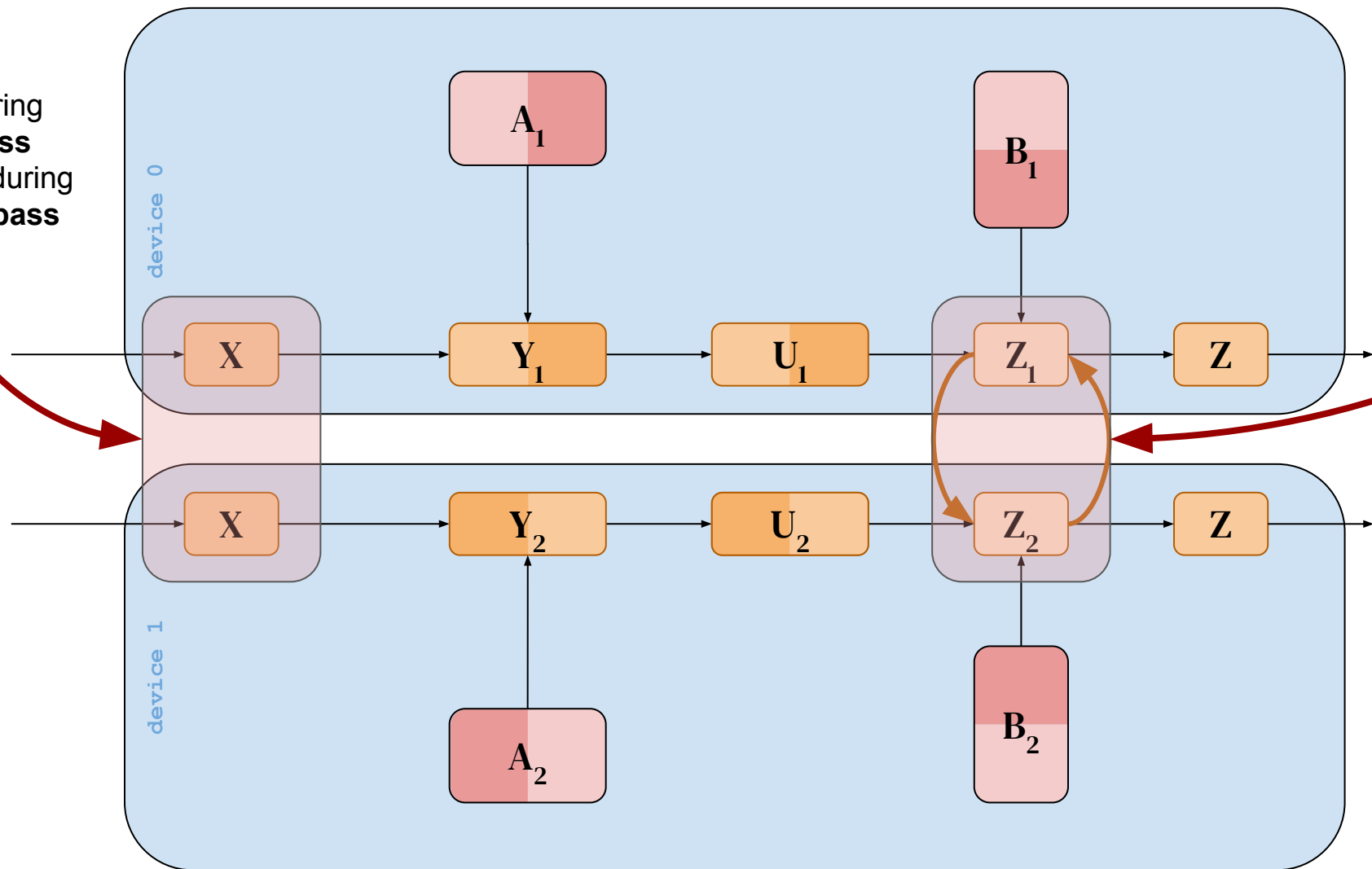




Forward pass

Painfully spelled out

- **Identity** during forward pass
- **AllReduce** during backward pass

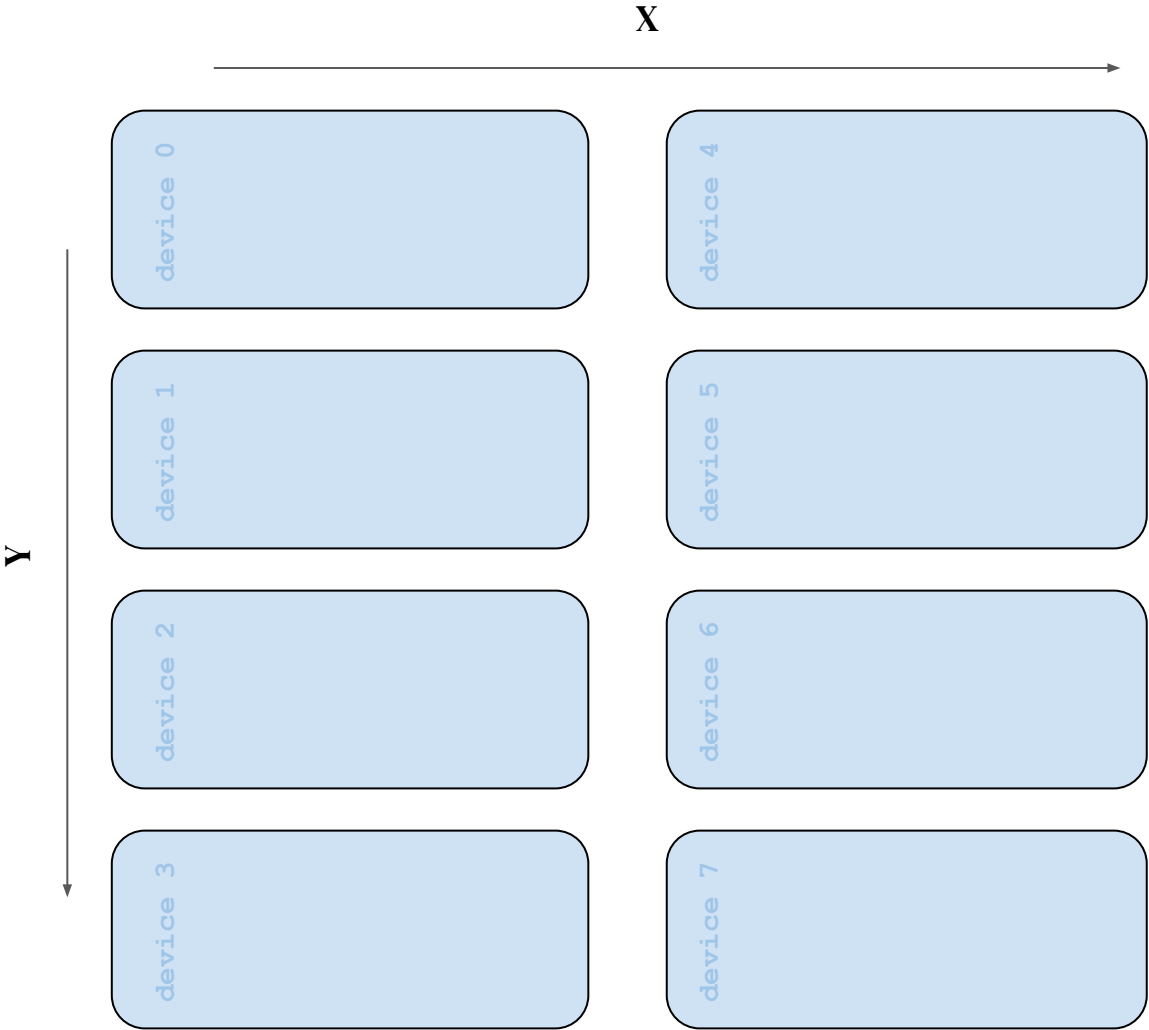
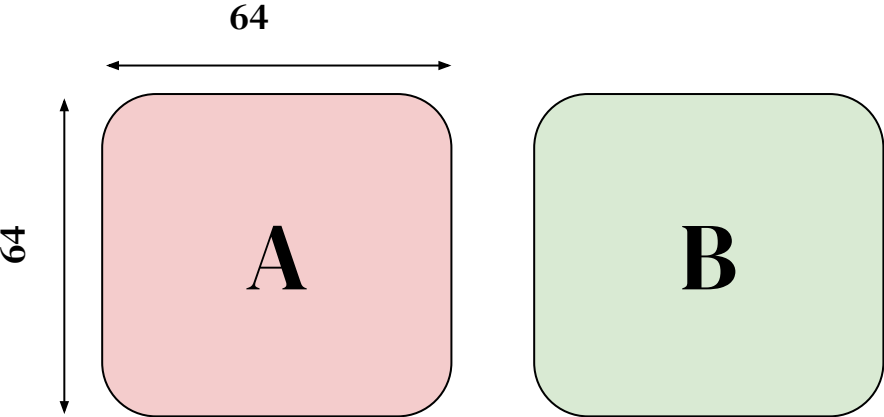


- **AllReduce** during forward pass
- **Identity** during backward pass

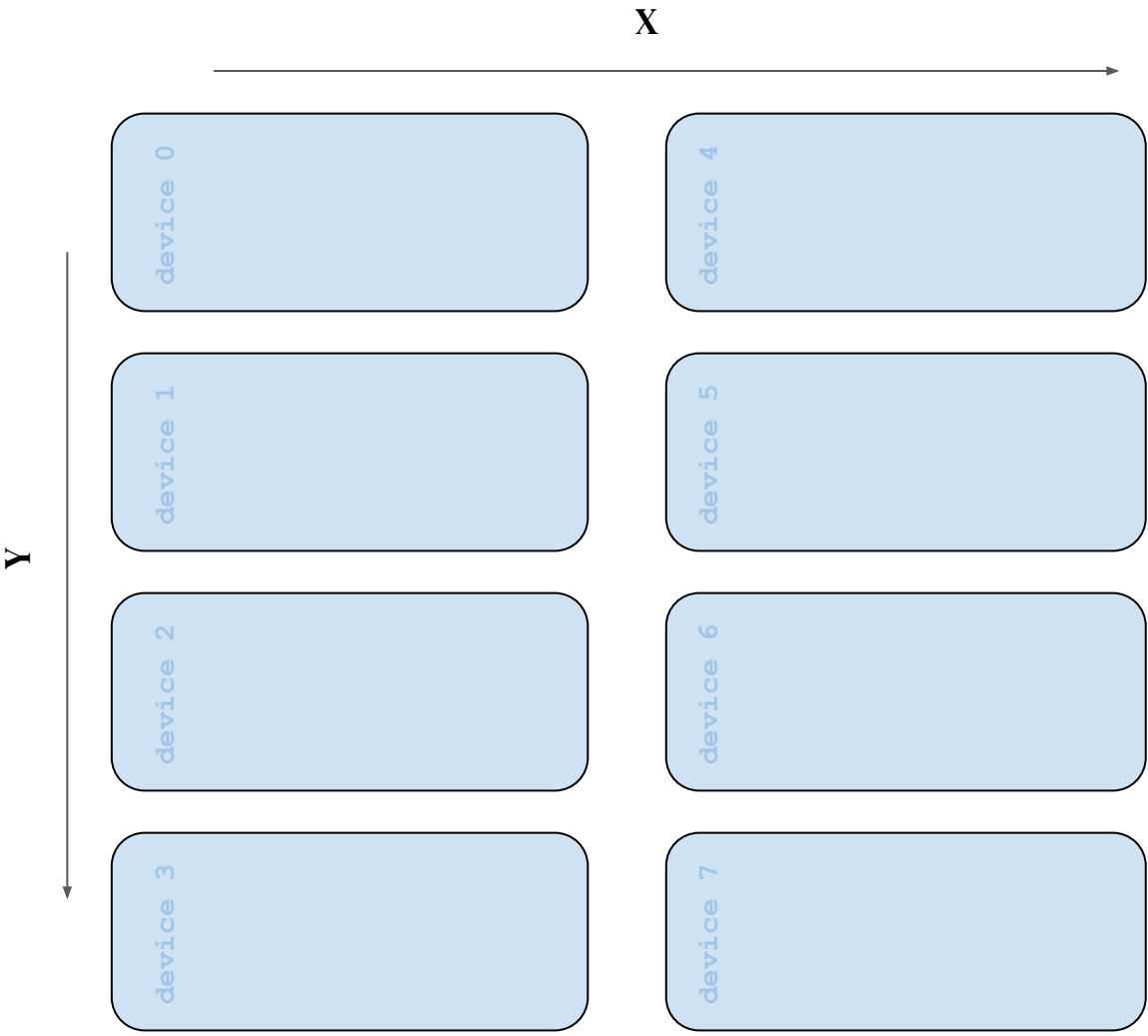
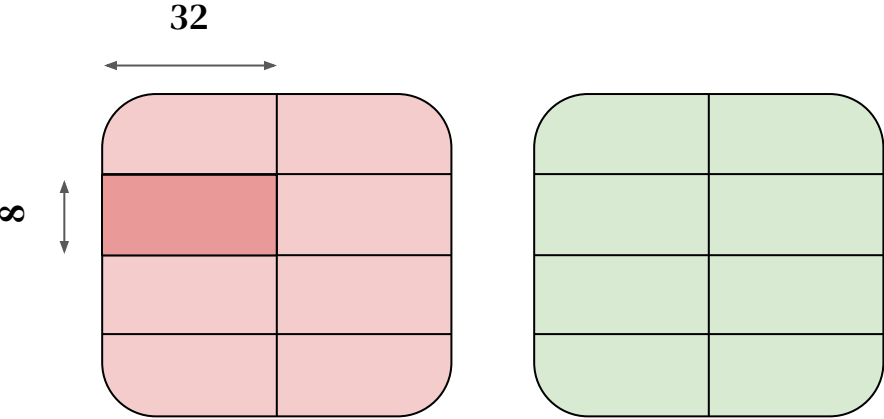
Sharded MLP - Summary

- Megatron sharding is a tensor parallelism technique, that avoids unnecessary communication between **N** devices.
- It can be applied to shard attention layers (splitting Q,K,V on heads).
- Reference: [Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism](#)

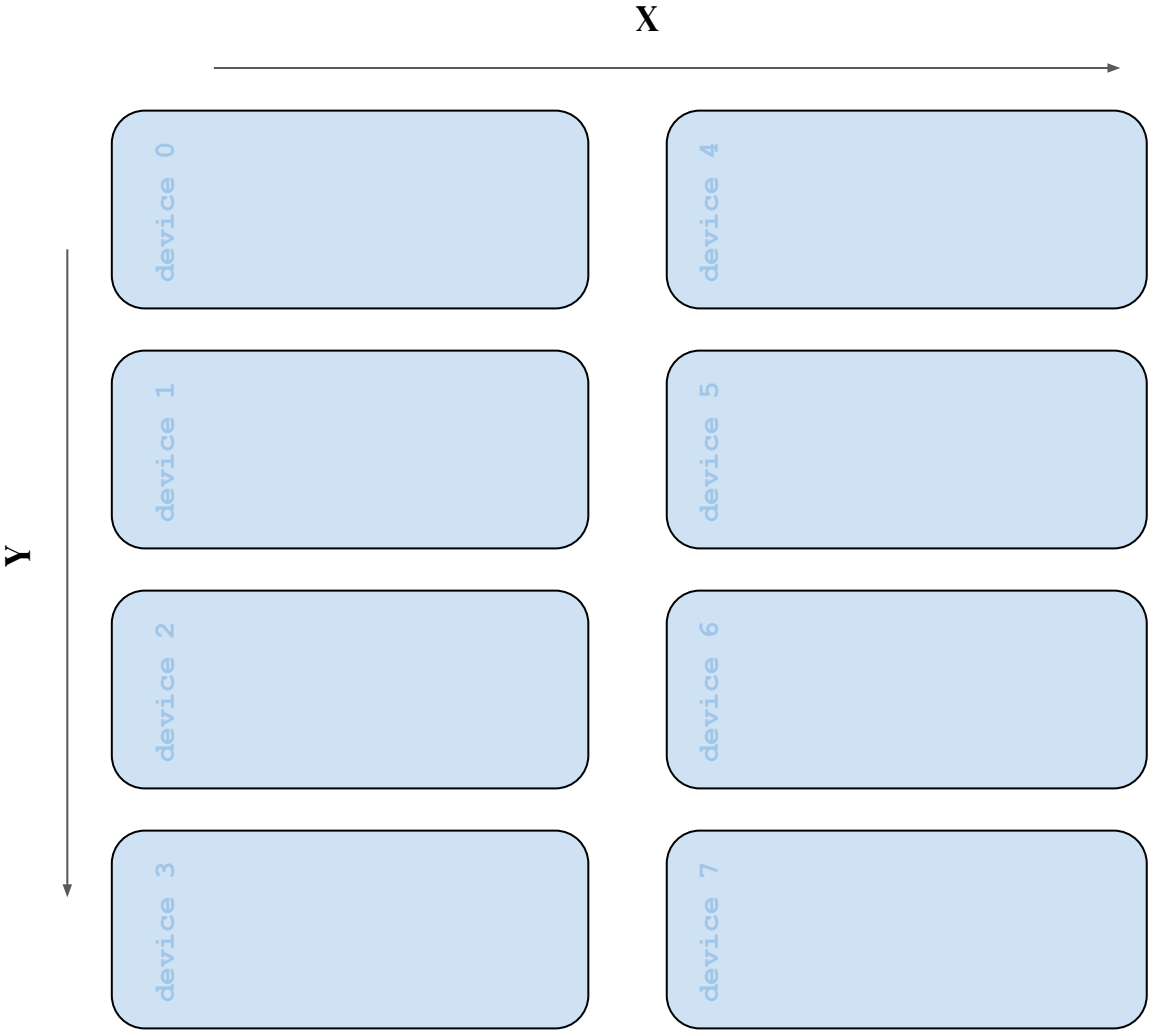
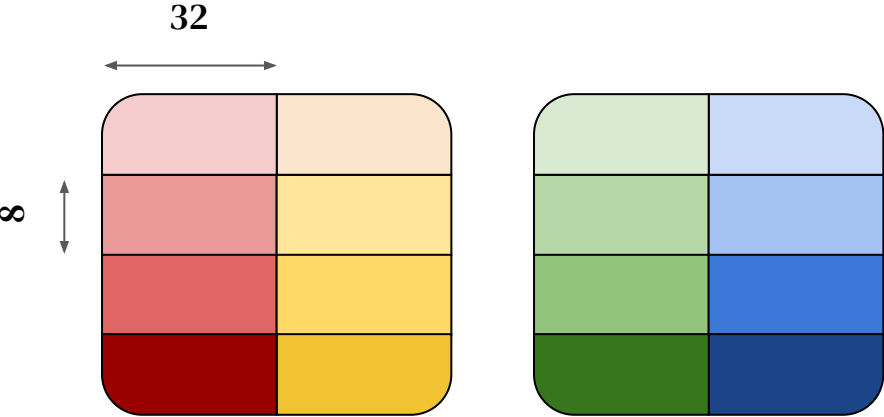
Matrix Multiplication on a mesh



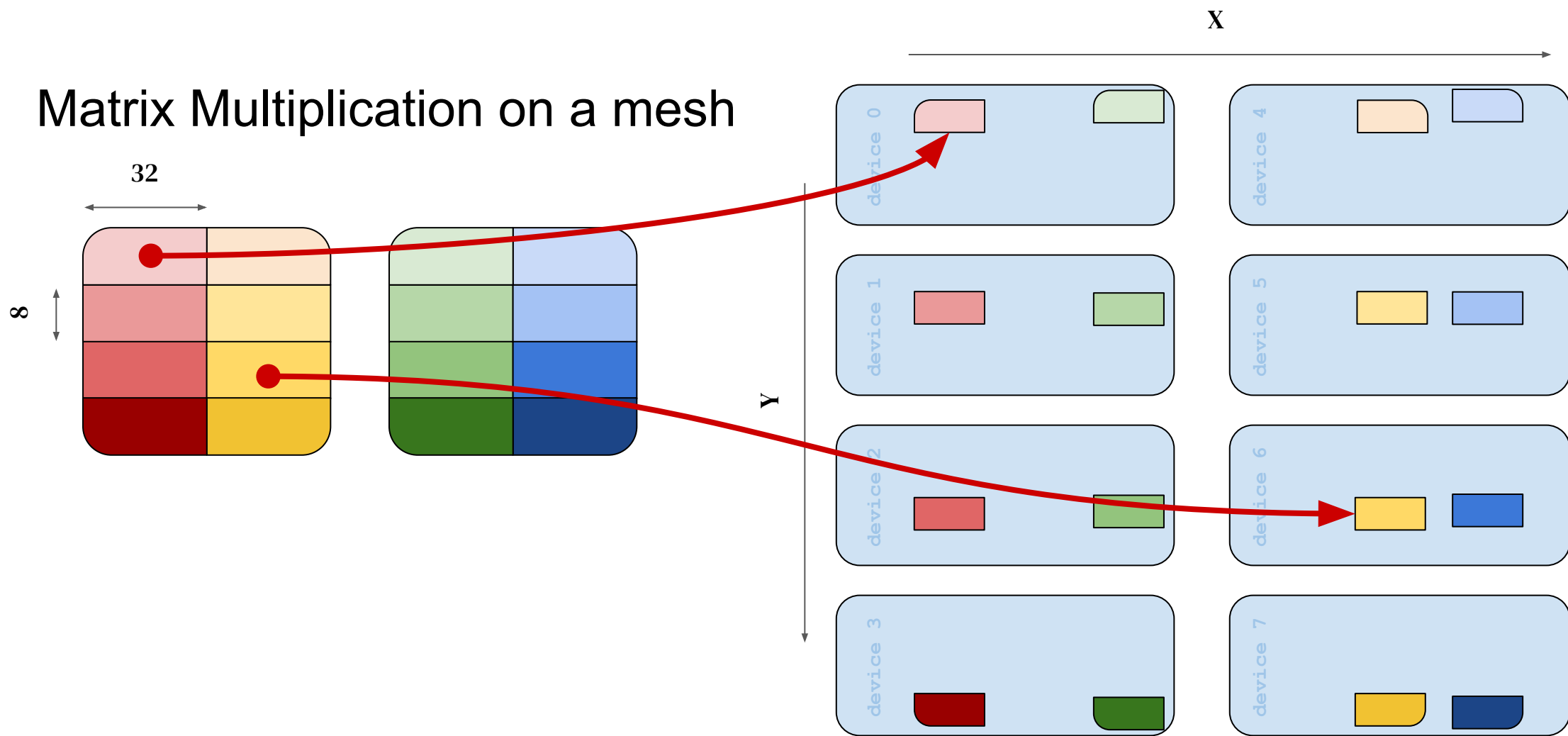
Matrix Multiplication on a mesh



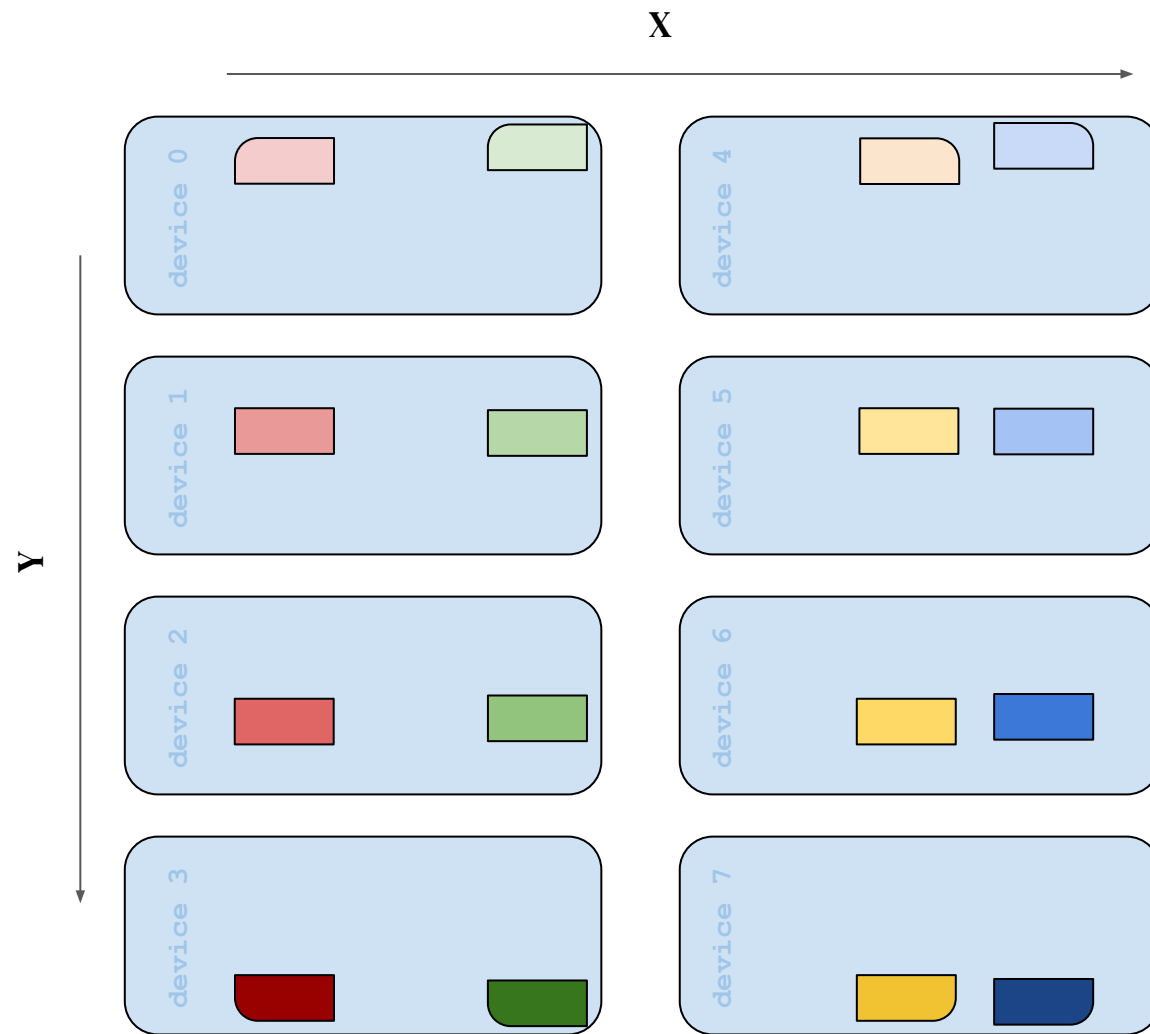
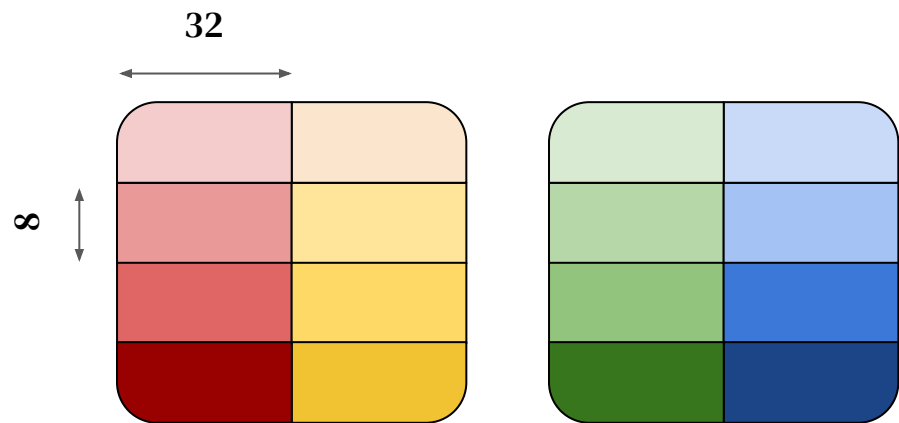
Matrix Multiplication on a mesh



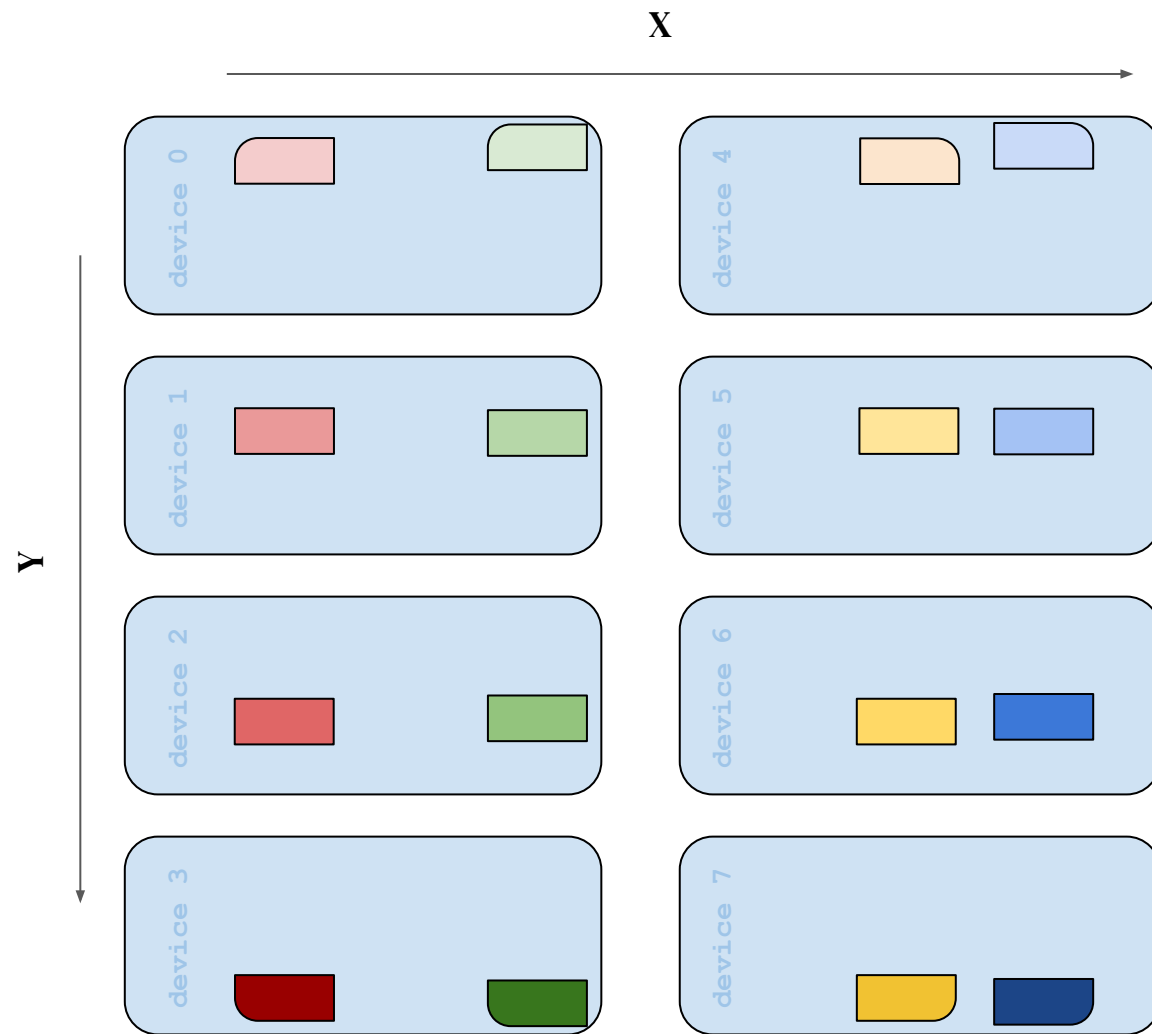
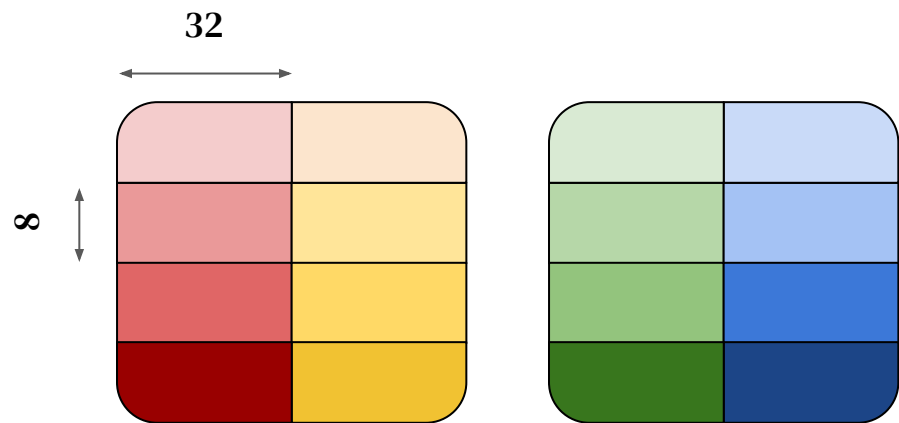
Matrix Multiplication on a mesh



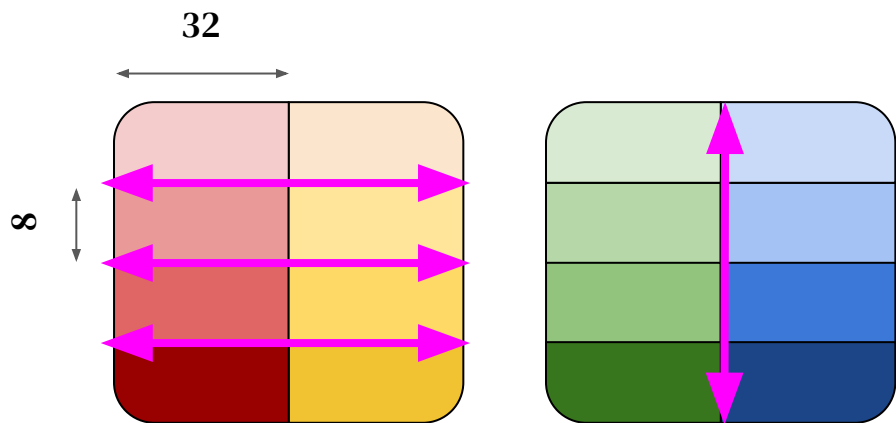
... but how do we multiply?



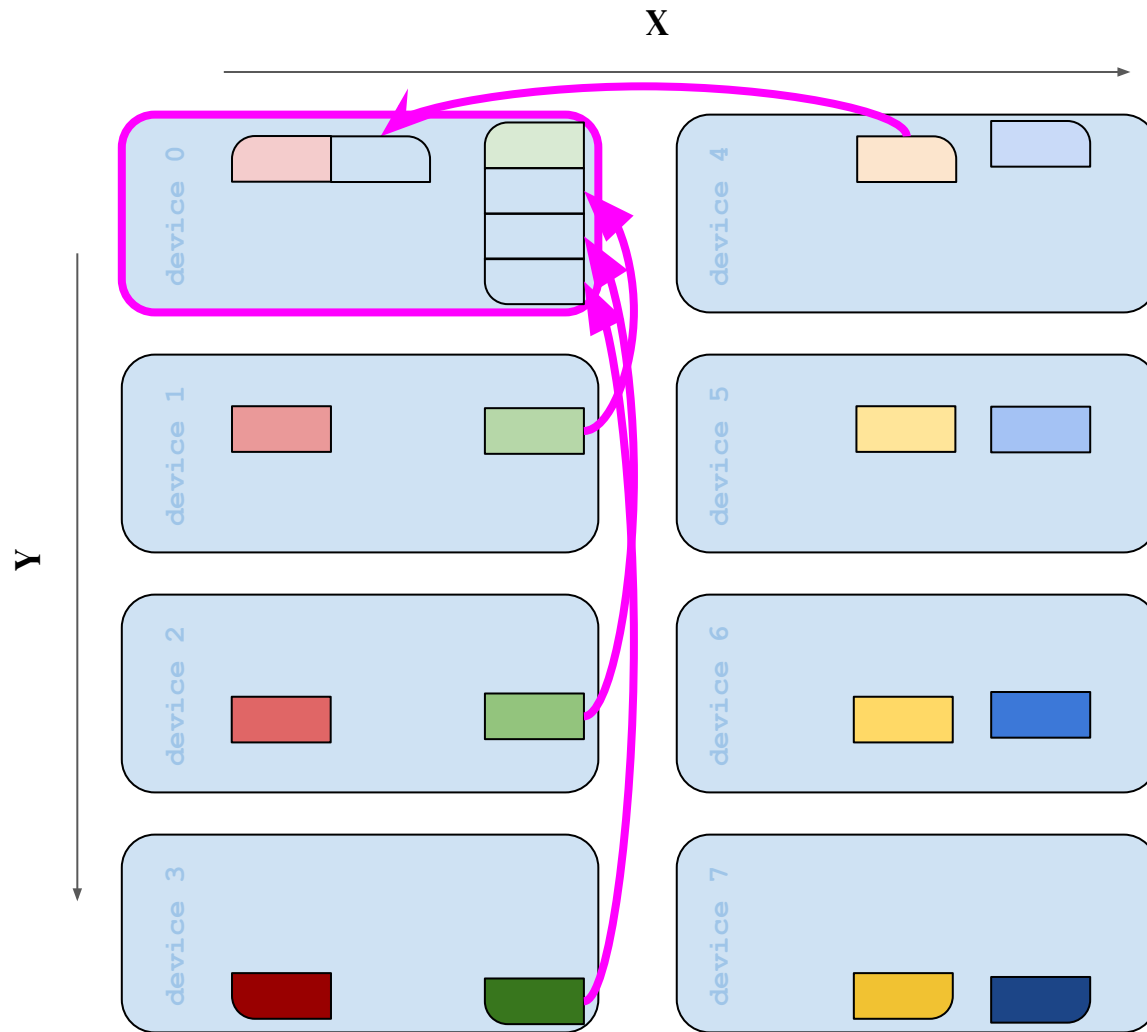
AllGather to the rescue



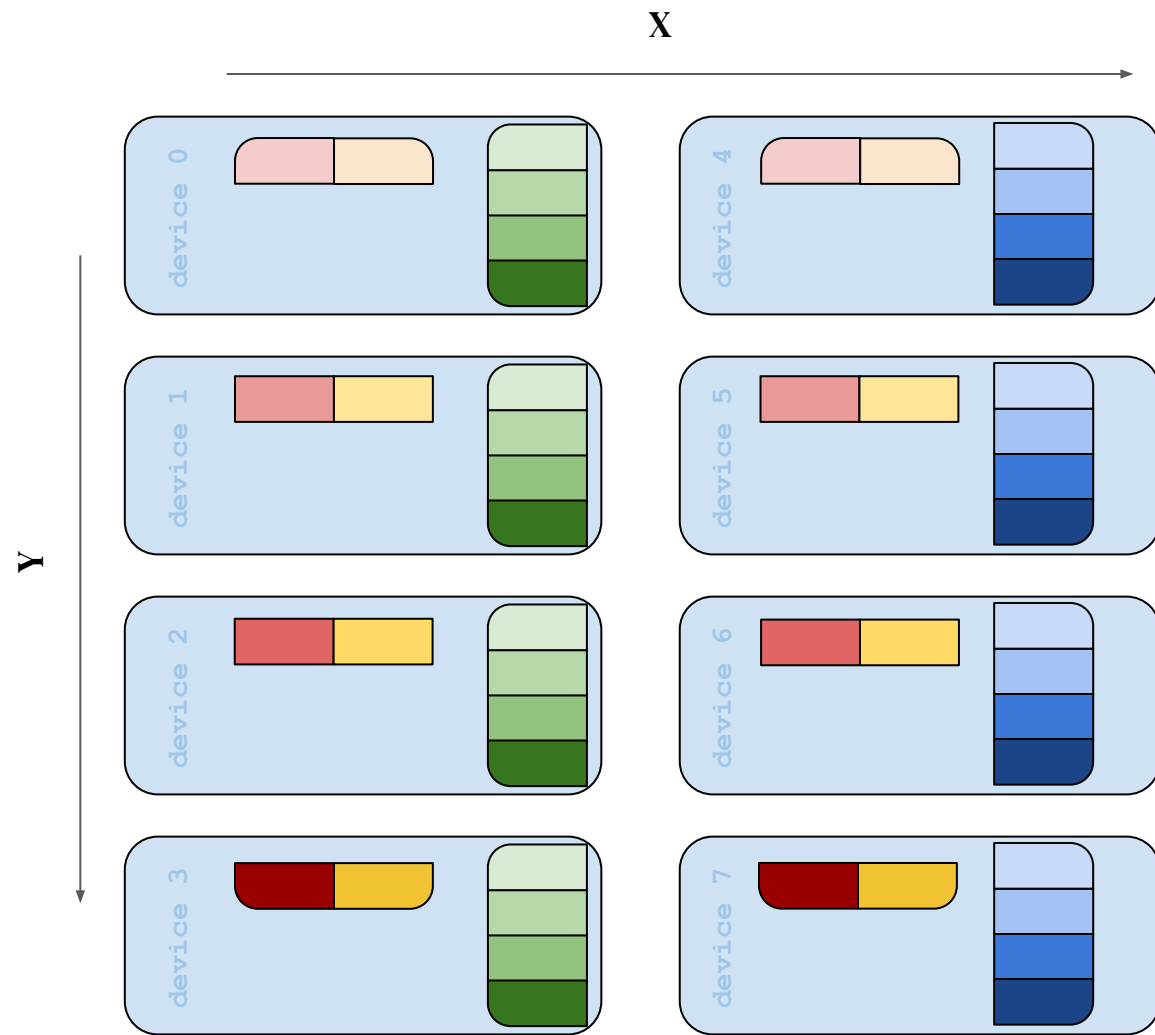
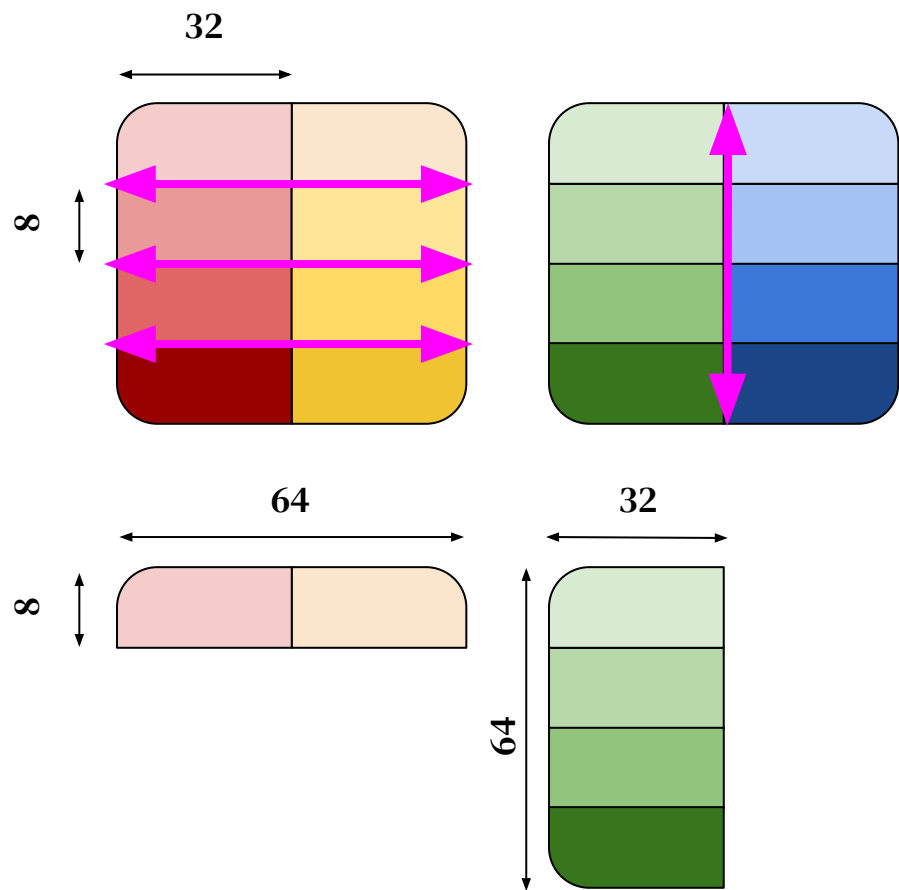
AllGather to the rescue



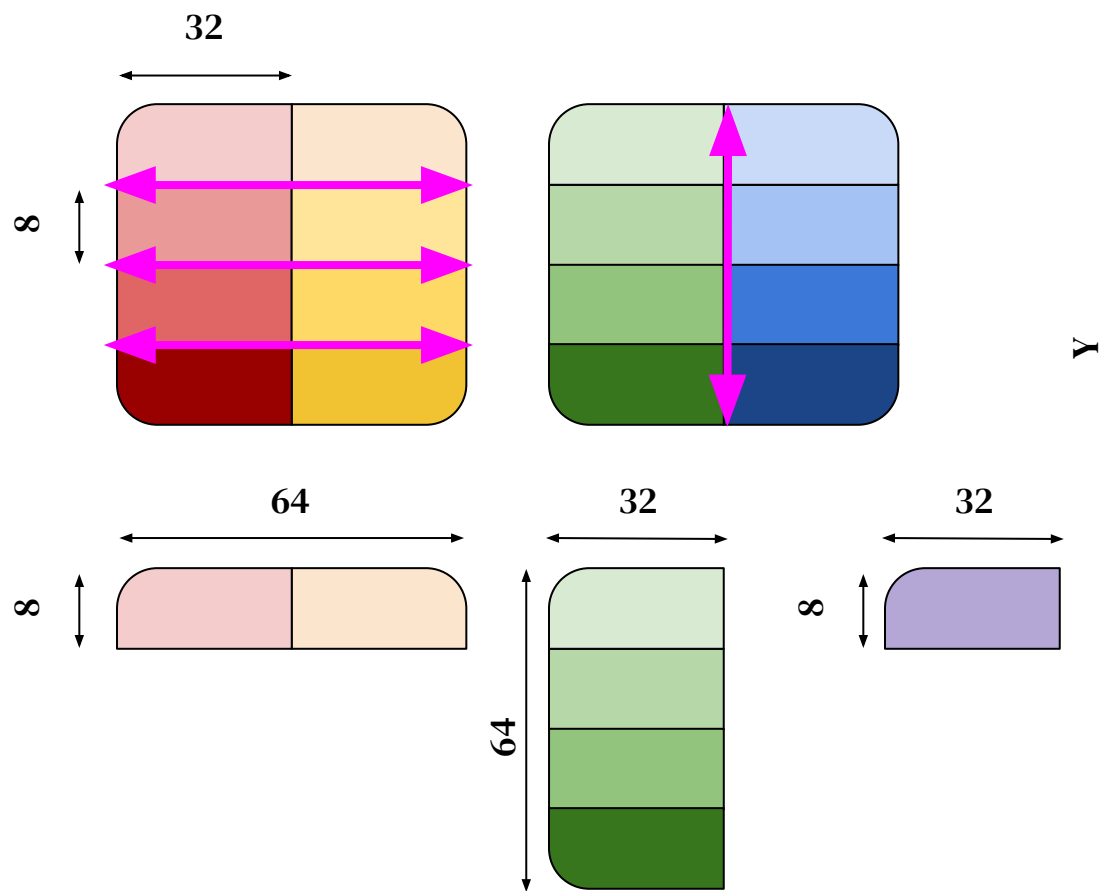
We're focusing on device 0, but this happens simultaneously on all devices.



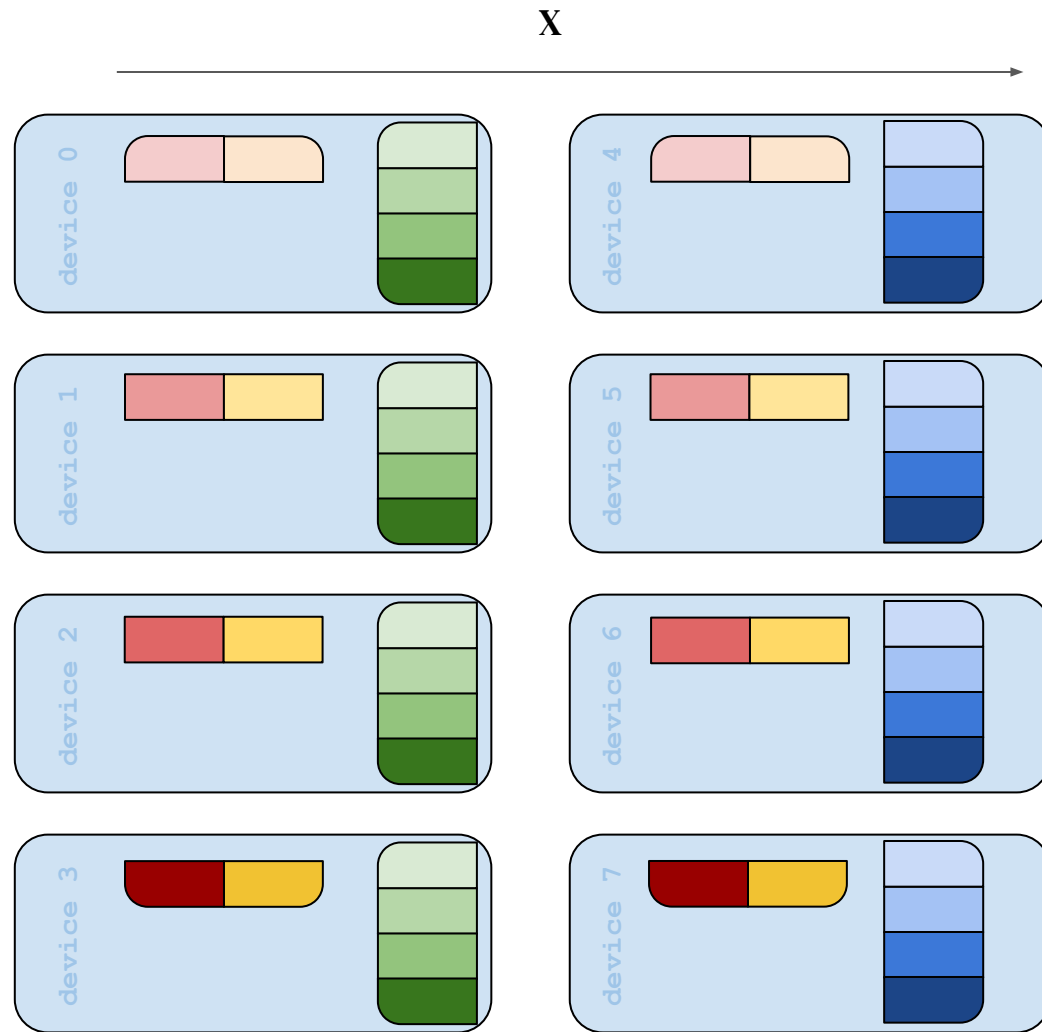
AllGather to the rescue



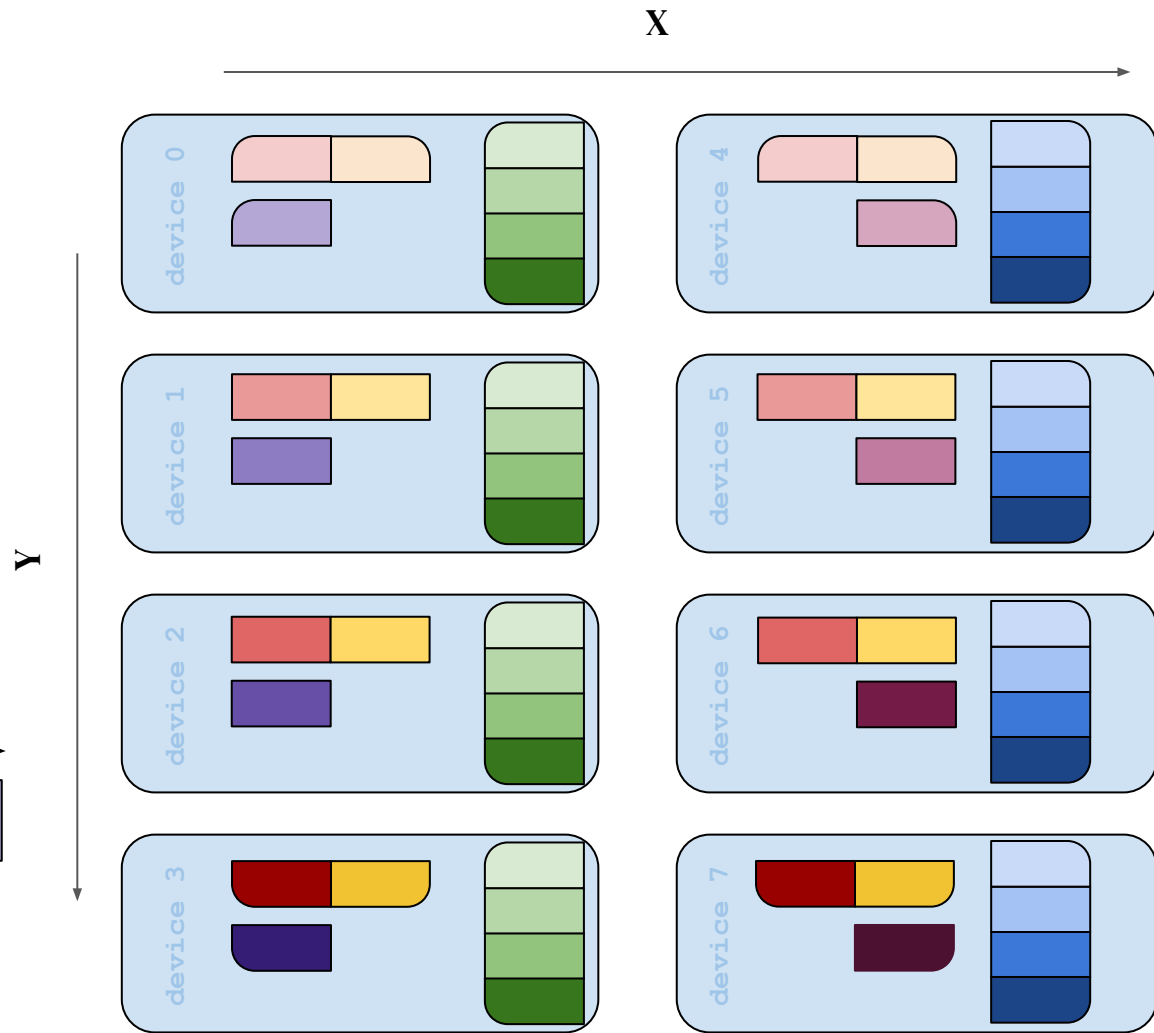
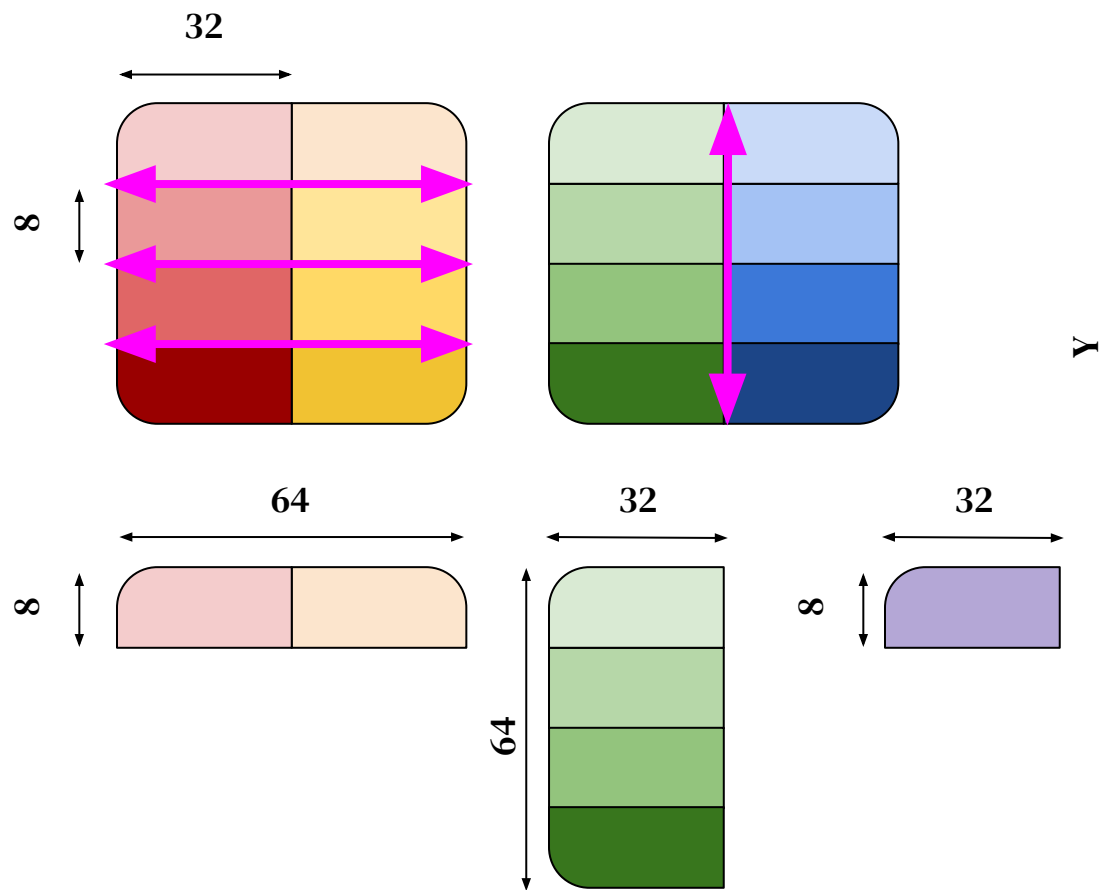
Now the shapes check out



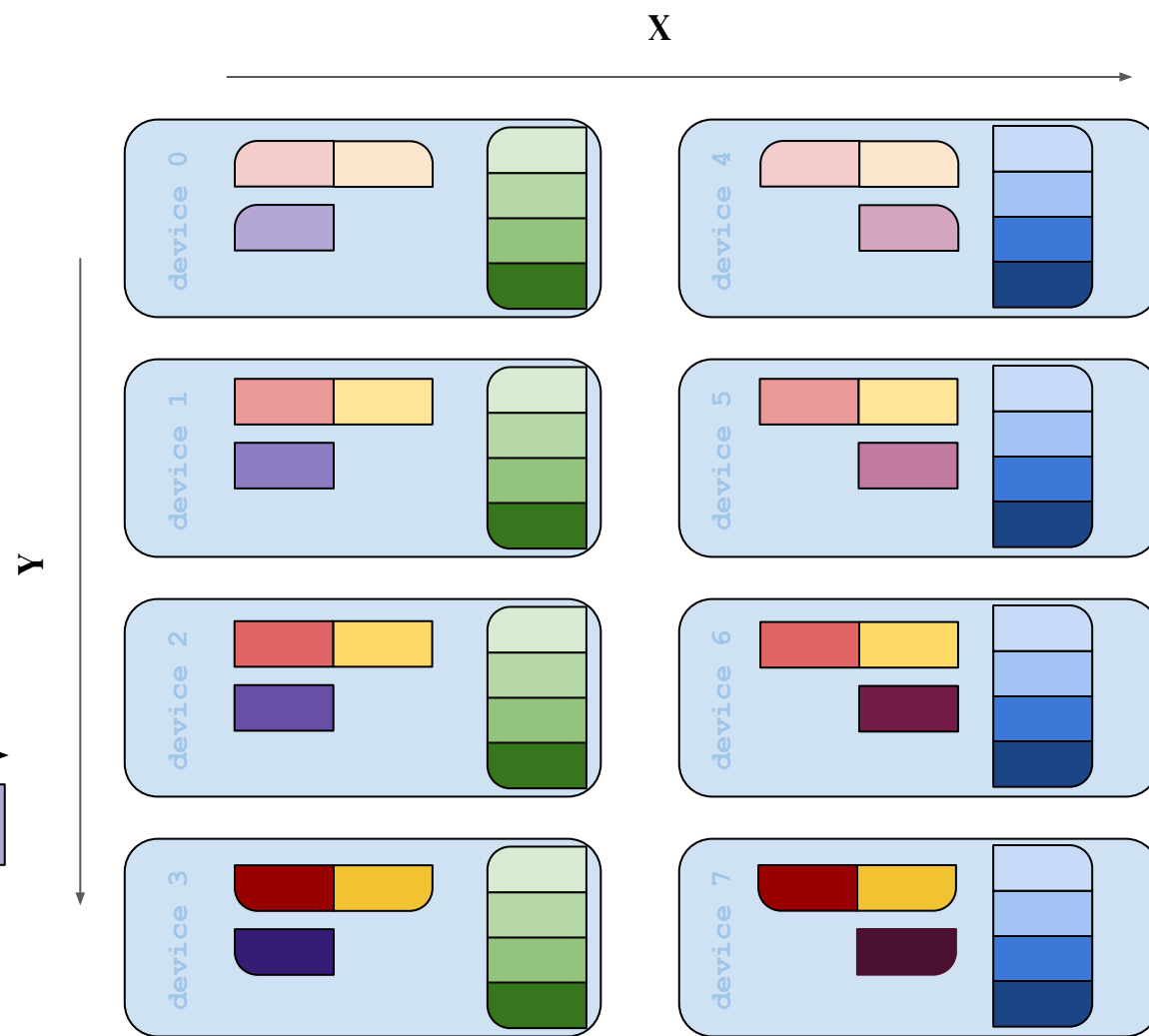
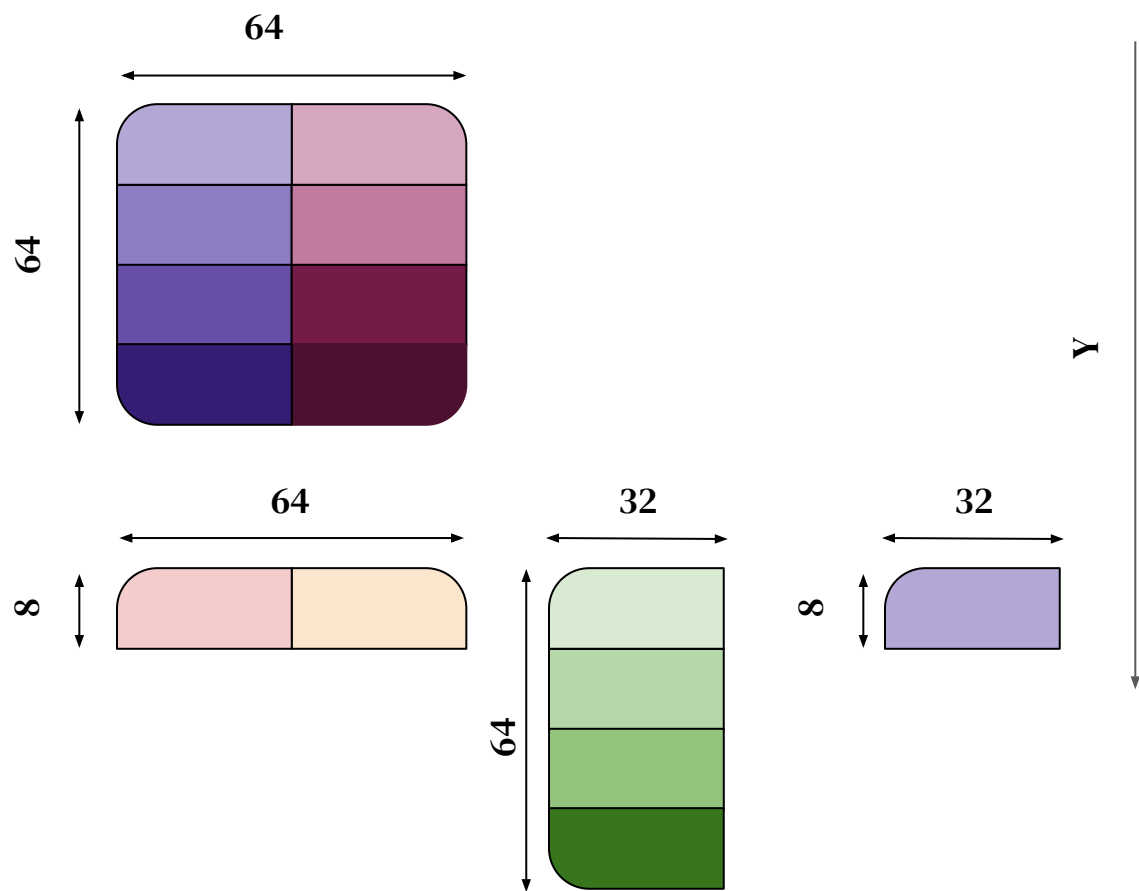
Y



Now the shapes check out

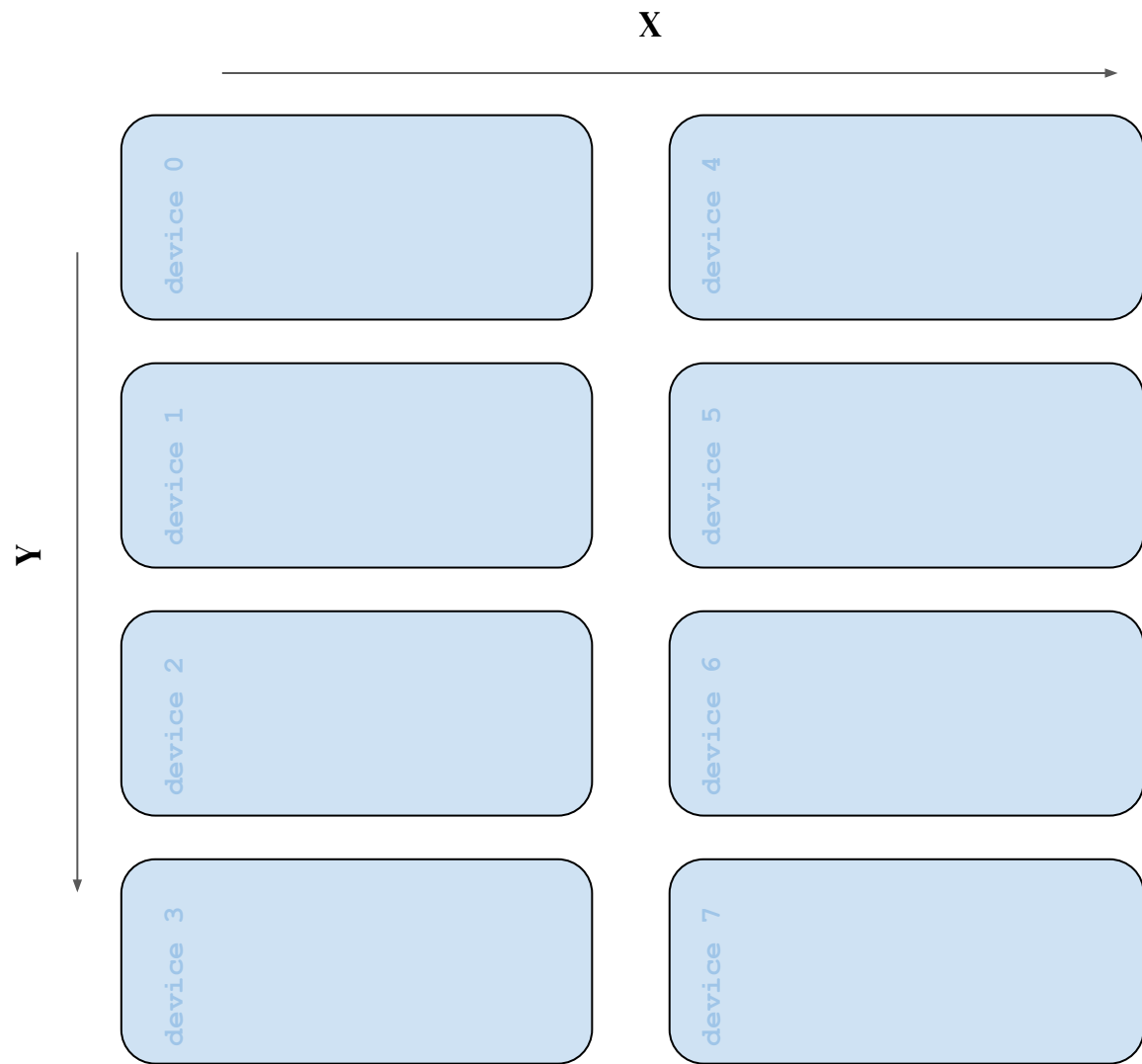
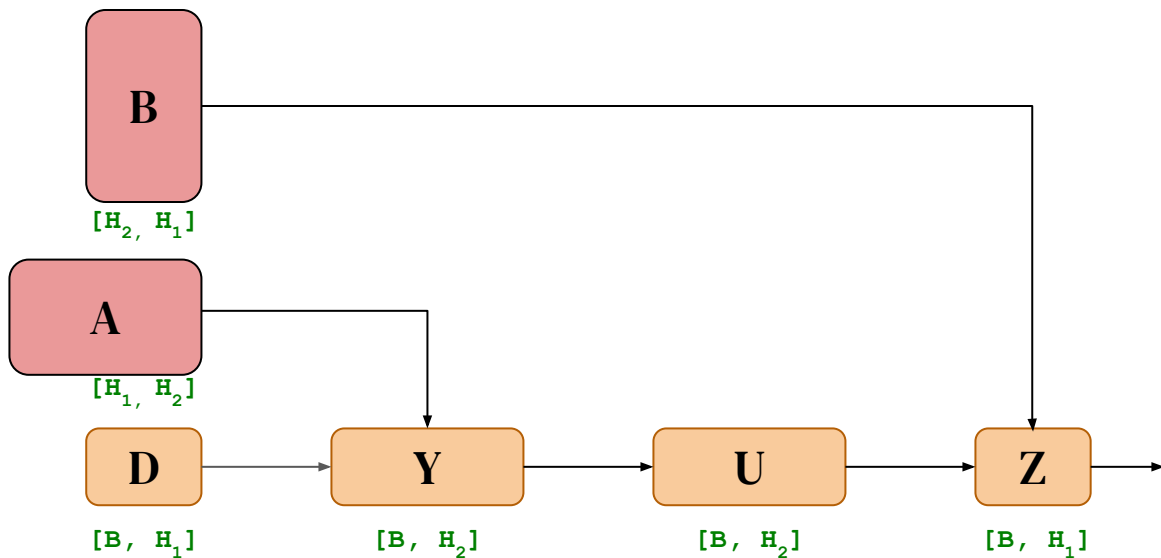


Result is also sharded

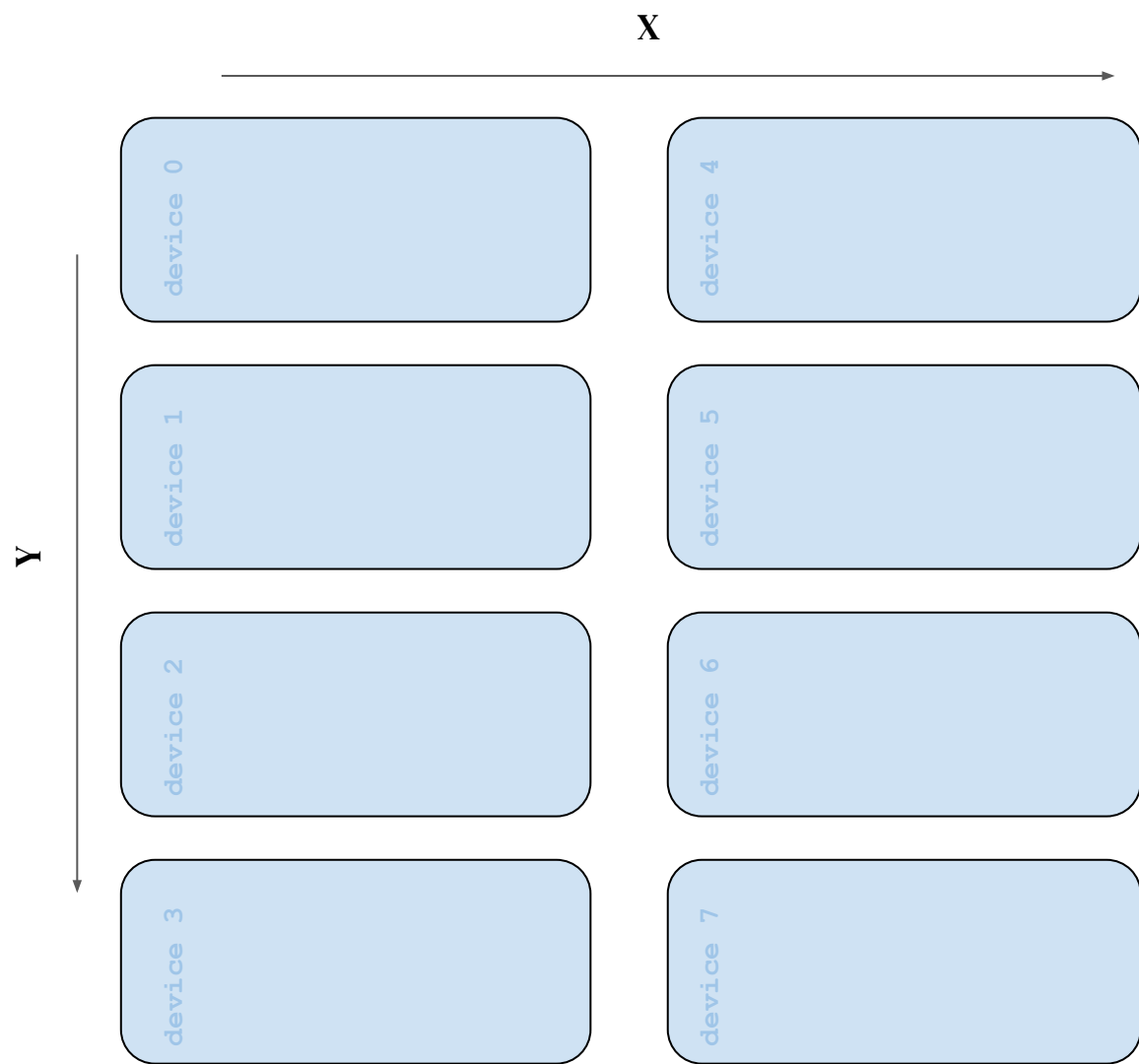
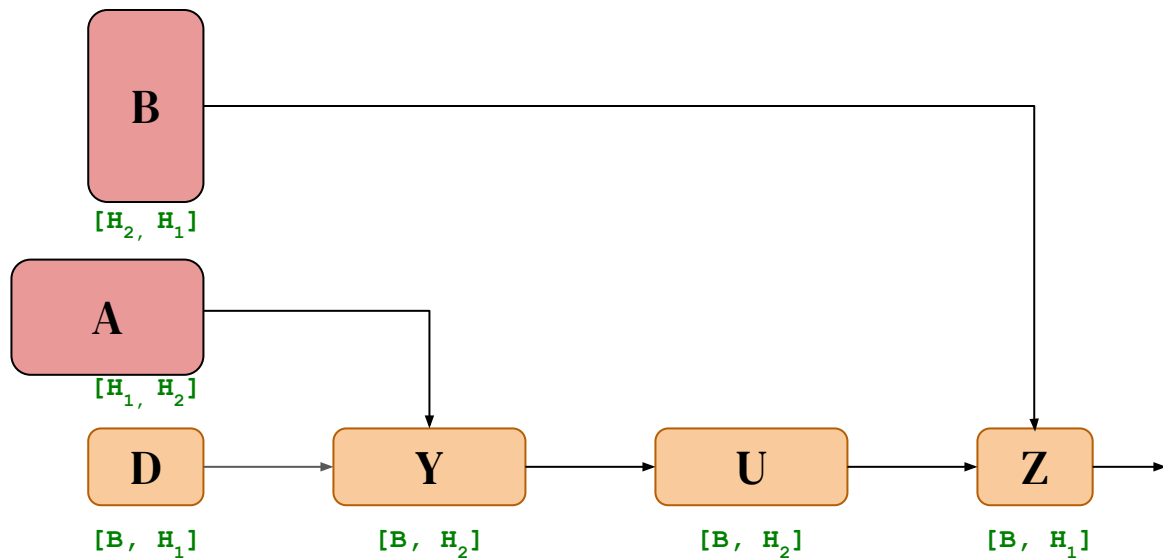


Sharding a 2-layer MLP

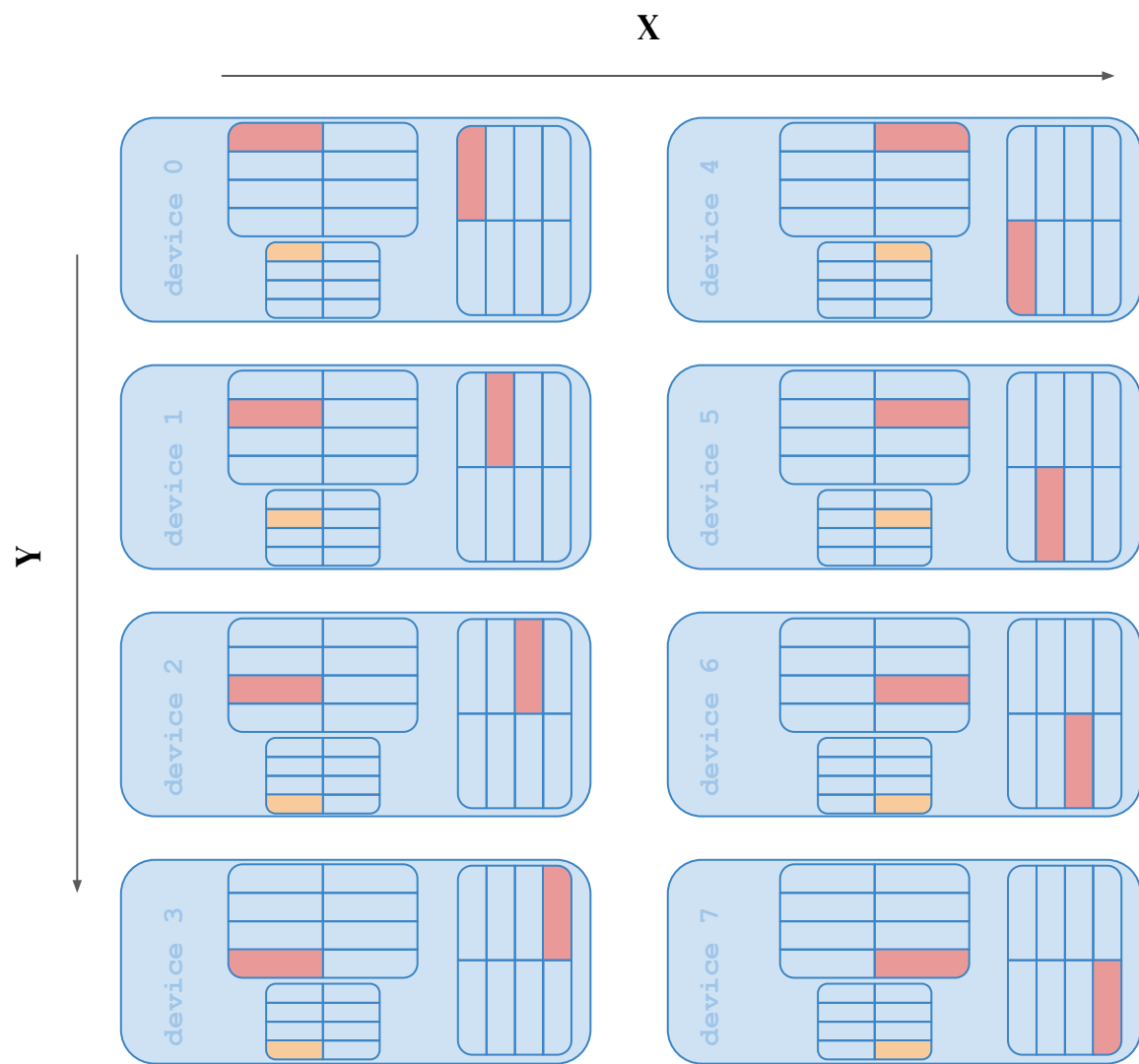
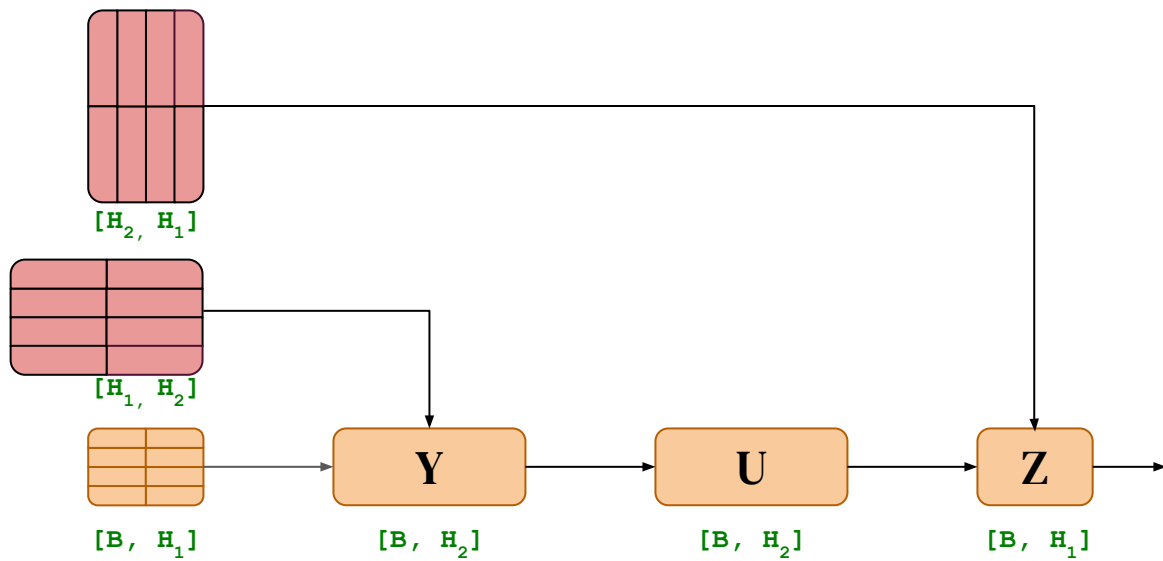
2D parallelism



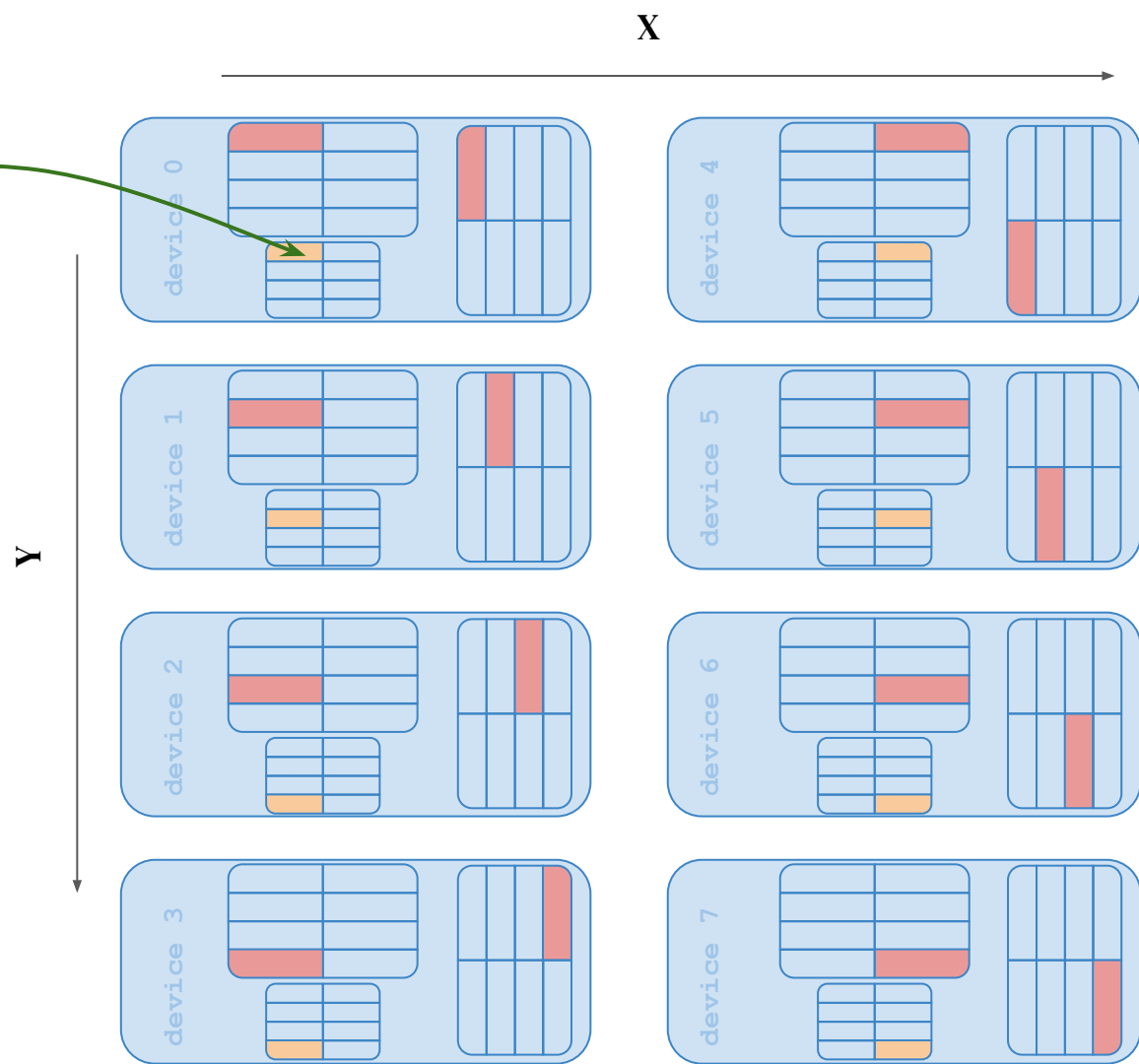
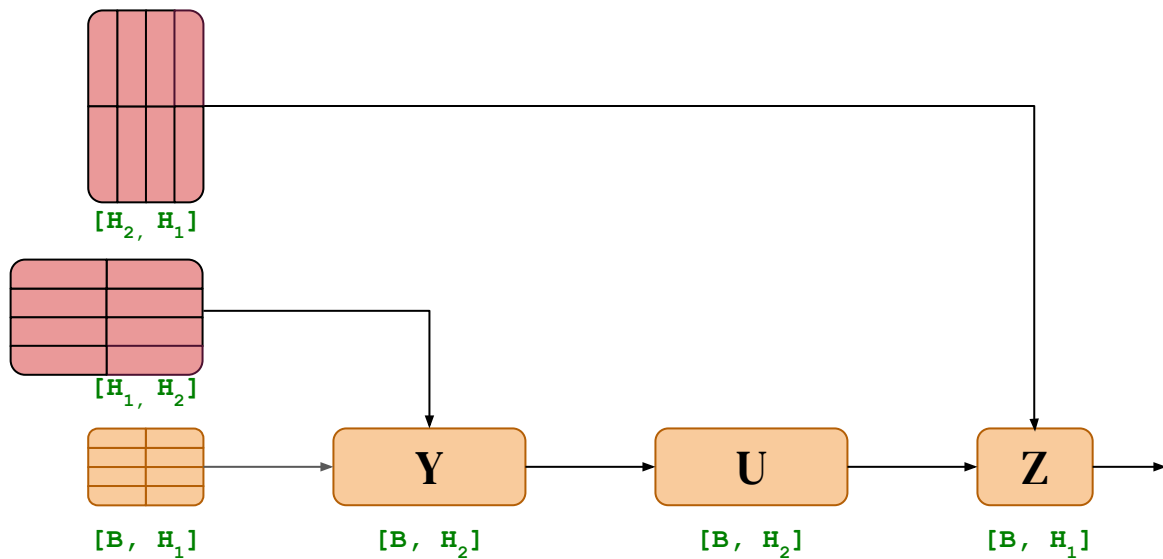
D_{YX} A_{YX} B_{XY}



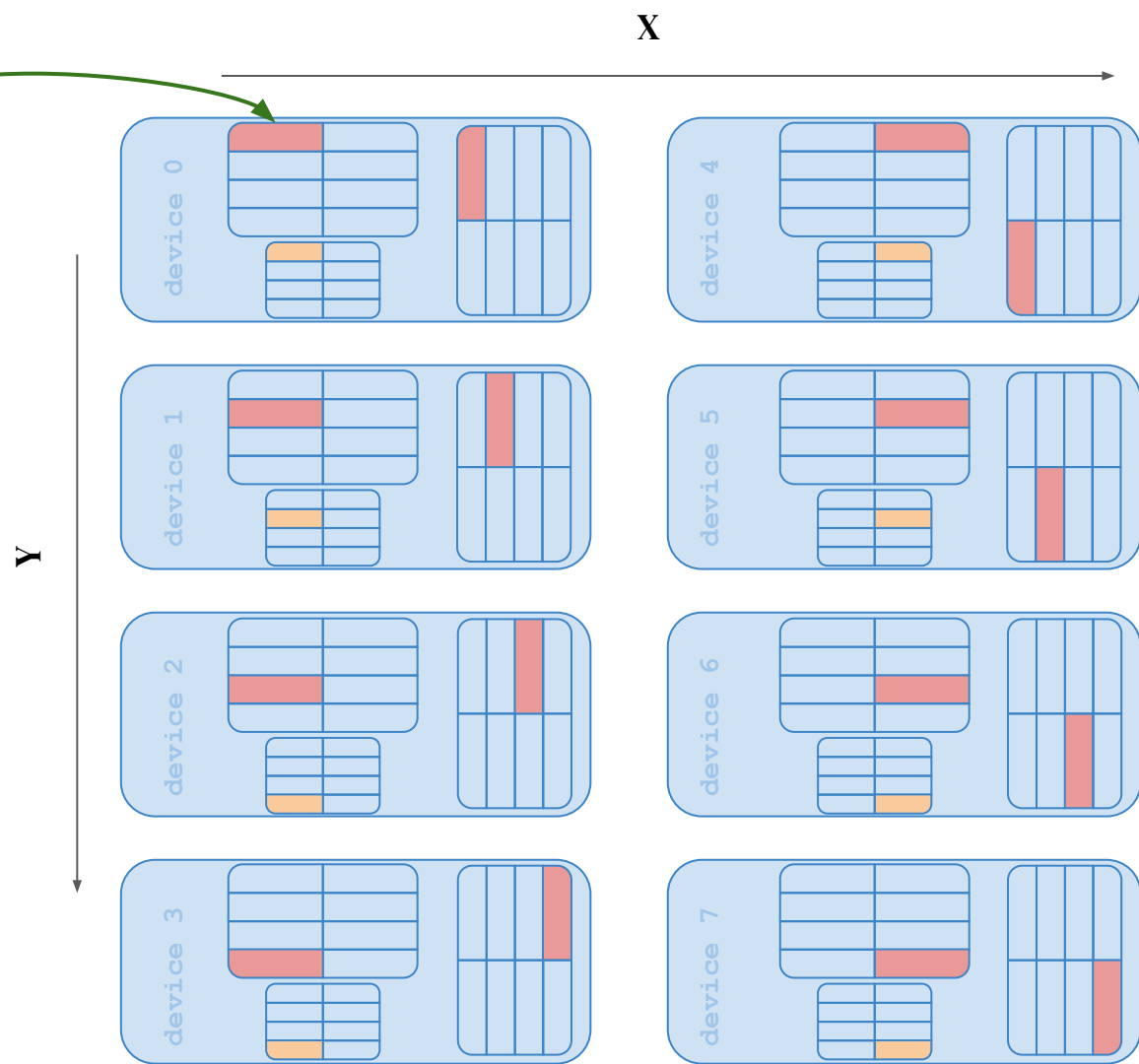
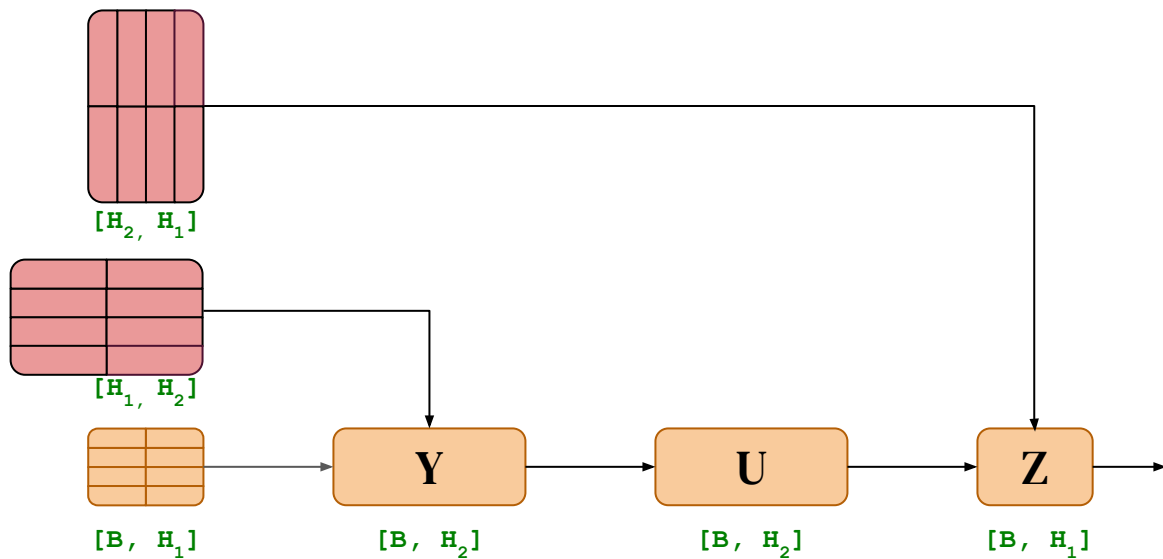
\mathbf{D}_{yx} \mathbf{A}_{yx} \mathbf{B}_{xy}



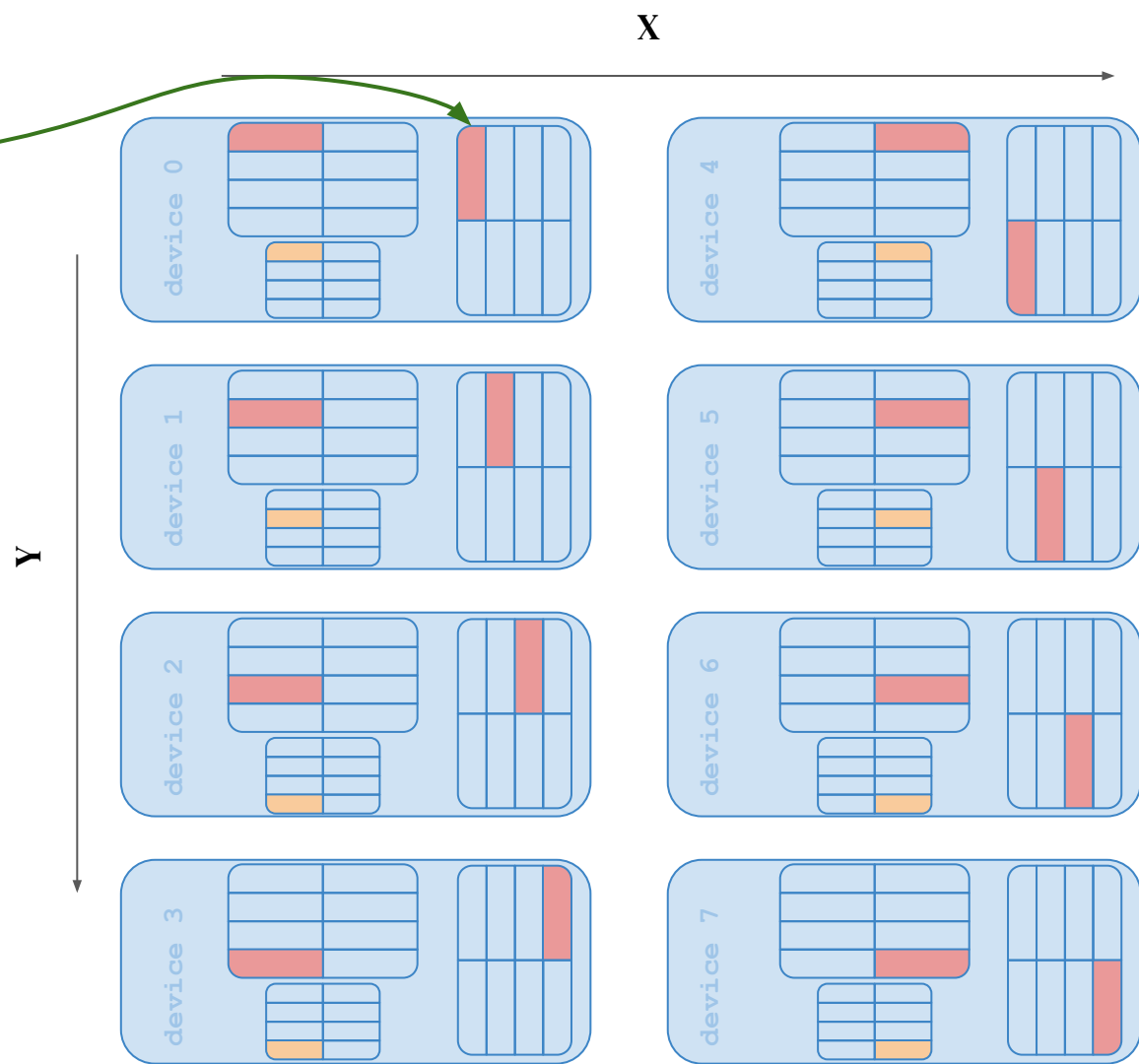
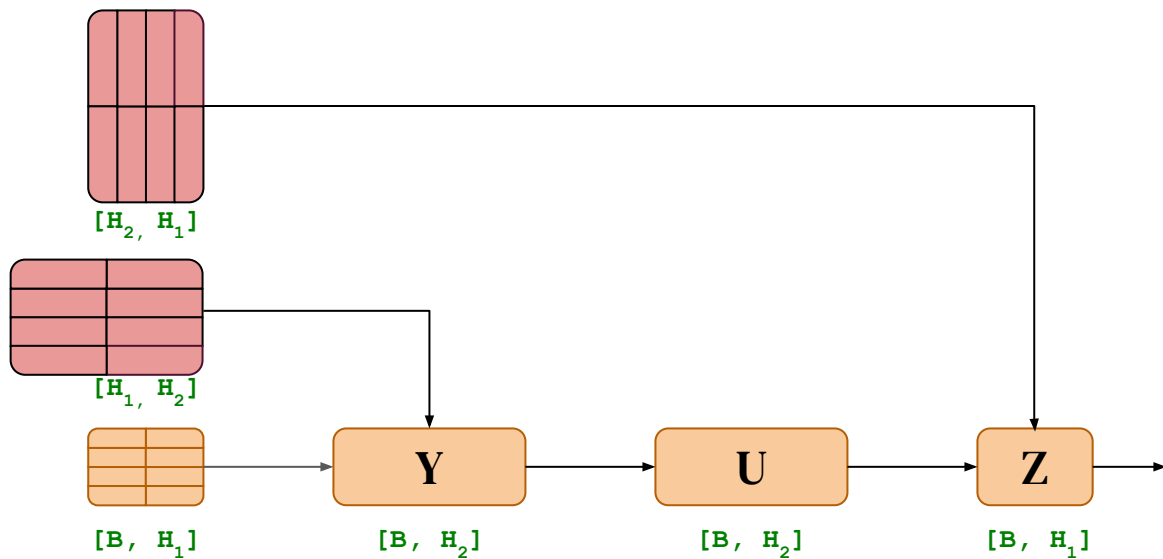
$$\boxed{D_{yx}} A_{yx} B_{xy} \quad [B/4, H_1/2]$$



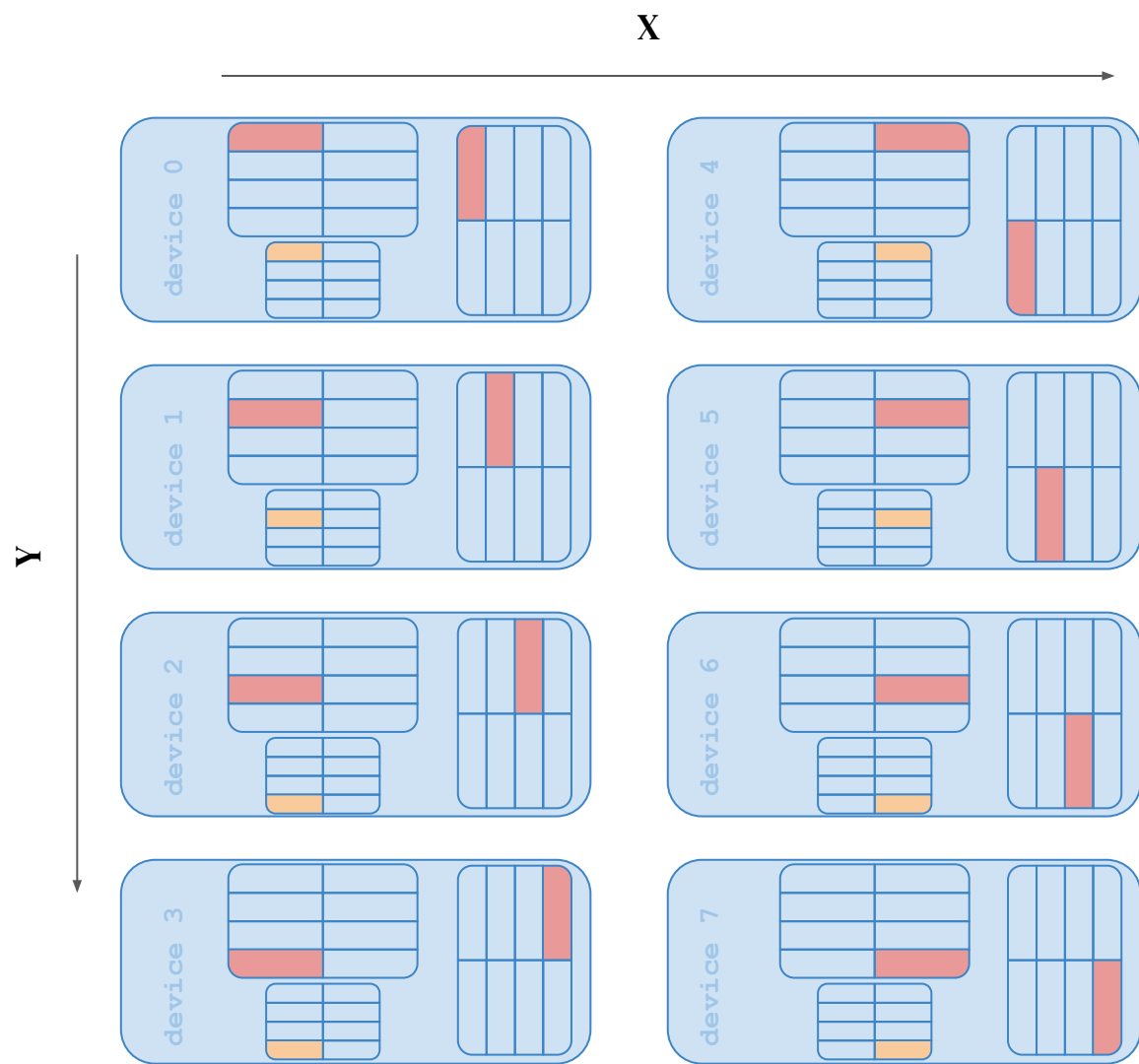
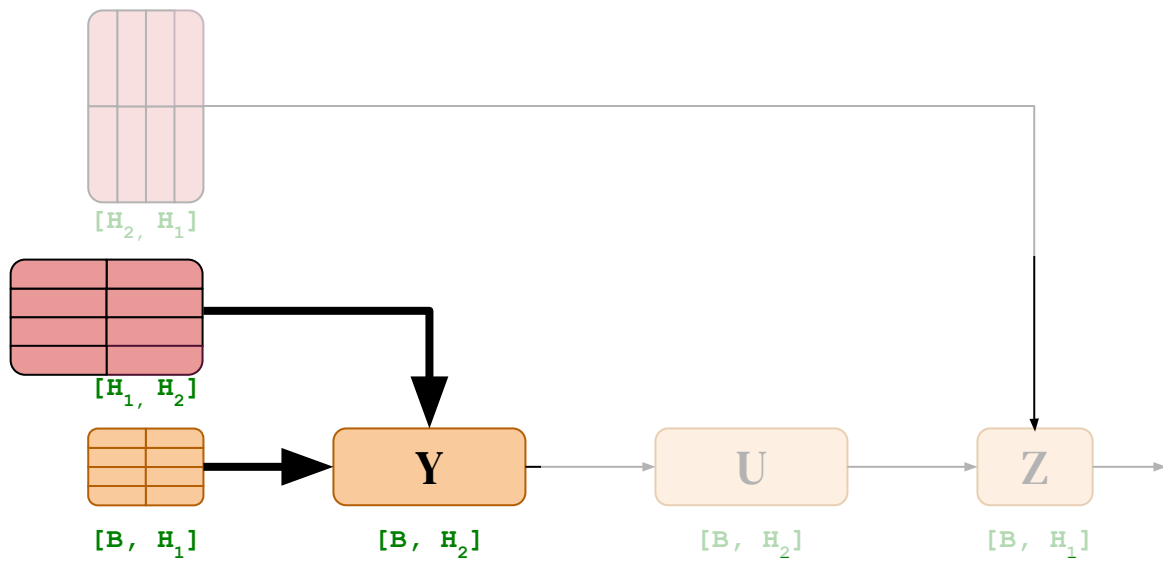
$$\mathbf{D}_{yx} \boxed{\mathbf{A}_{yx}} \mathbf{B}_{xy} \quad [\mathbf{H}_1/4, \mathbf{H}_2/2]$$



\mathbf{D}_{yx} \mathbf{A}_{yx} \mathbf{B}_{xy} $[\mathbf{H}_2/2, \mathbf{H}_1/4]$

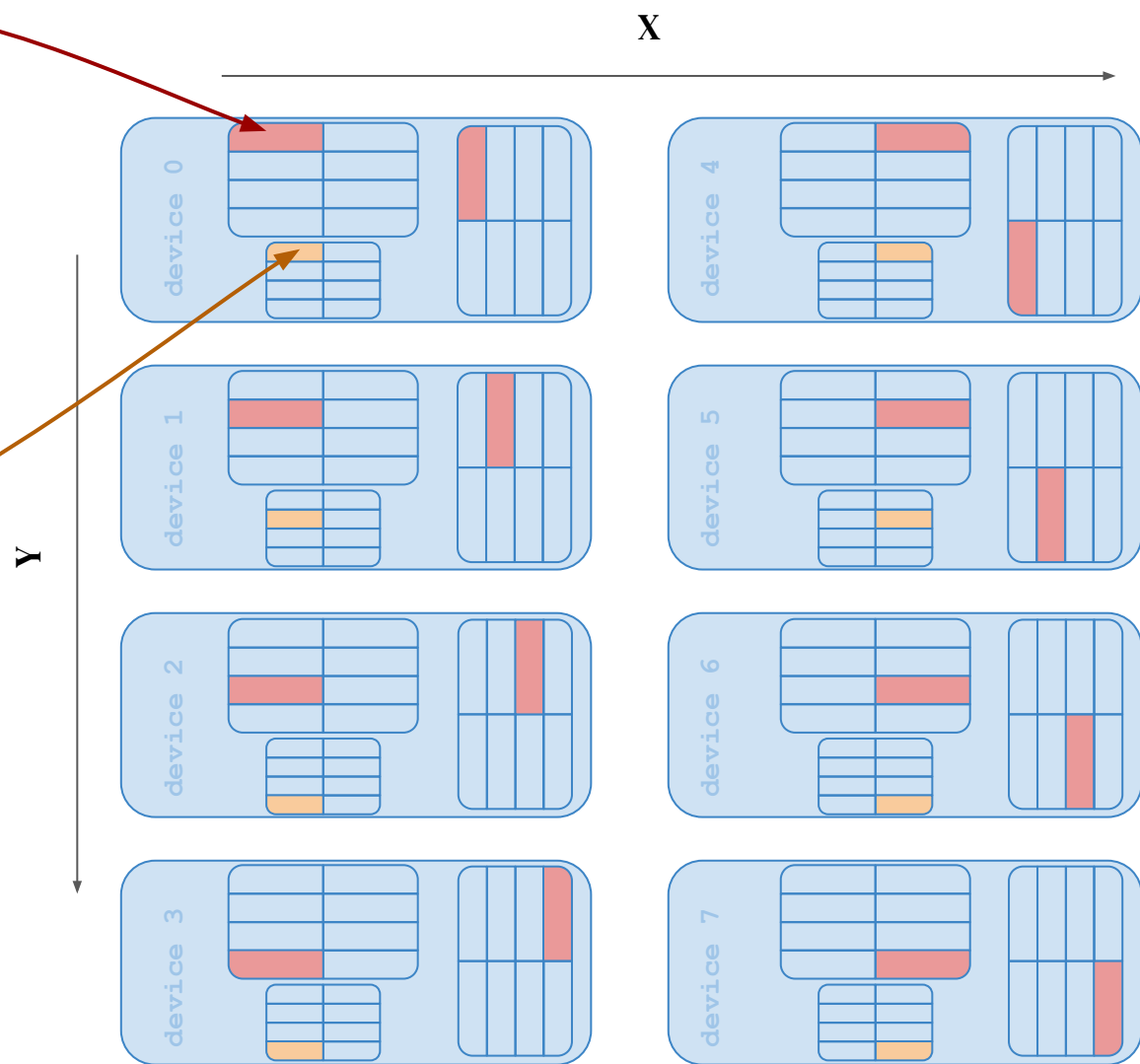
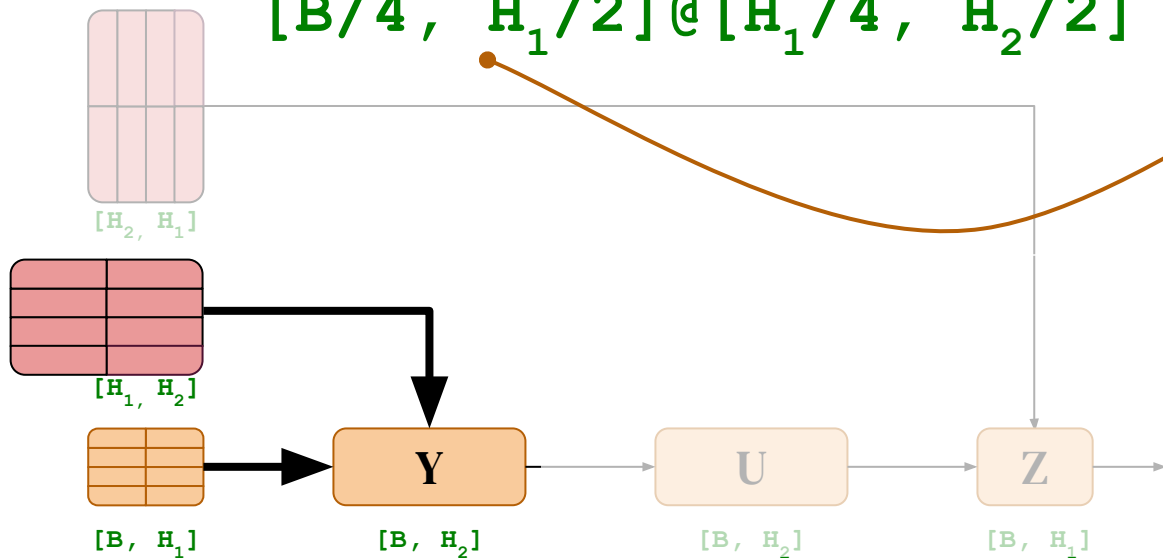


Compute $\mathbf{D}_{yx} @ \mathbf{A}_{yx}$

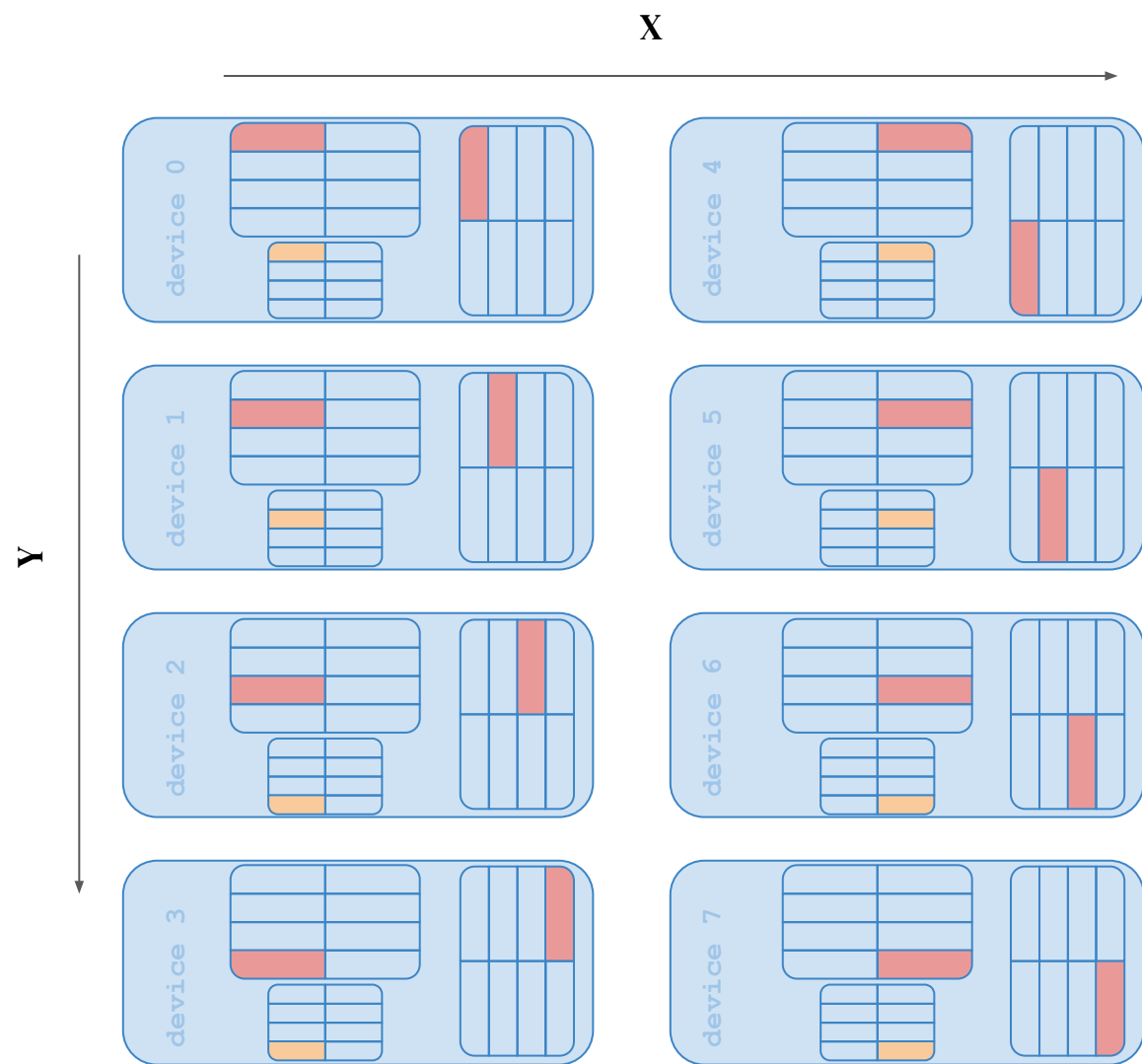
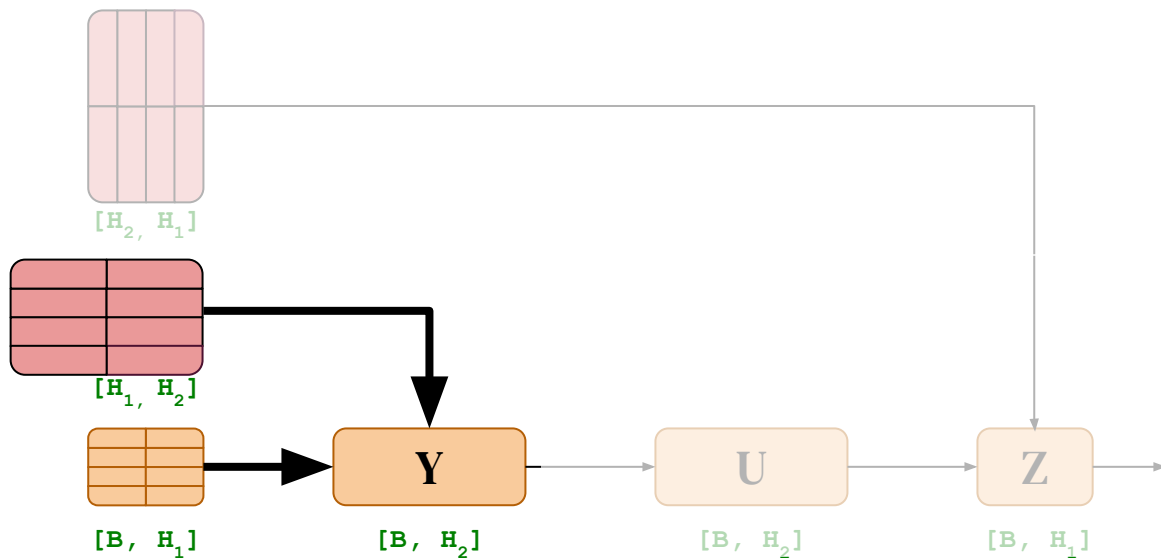


Compute \mathbf{D}_{yx} @ \mathbf{A}_{yx}

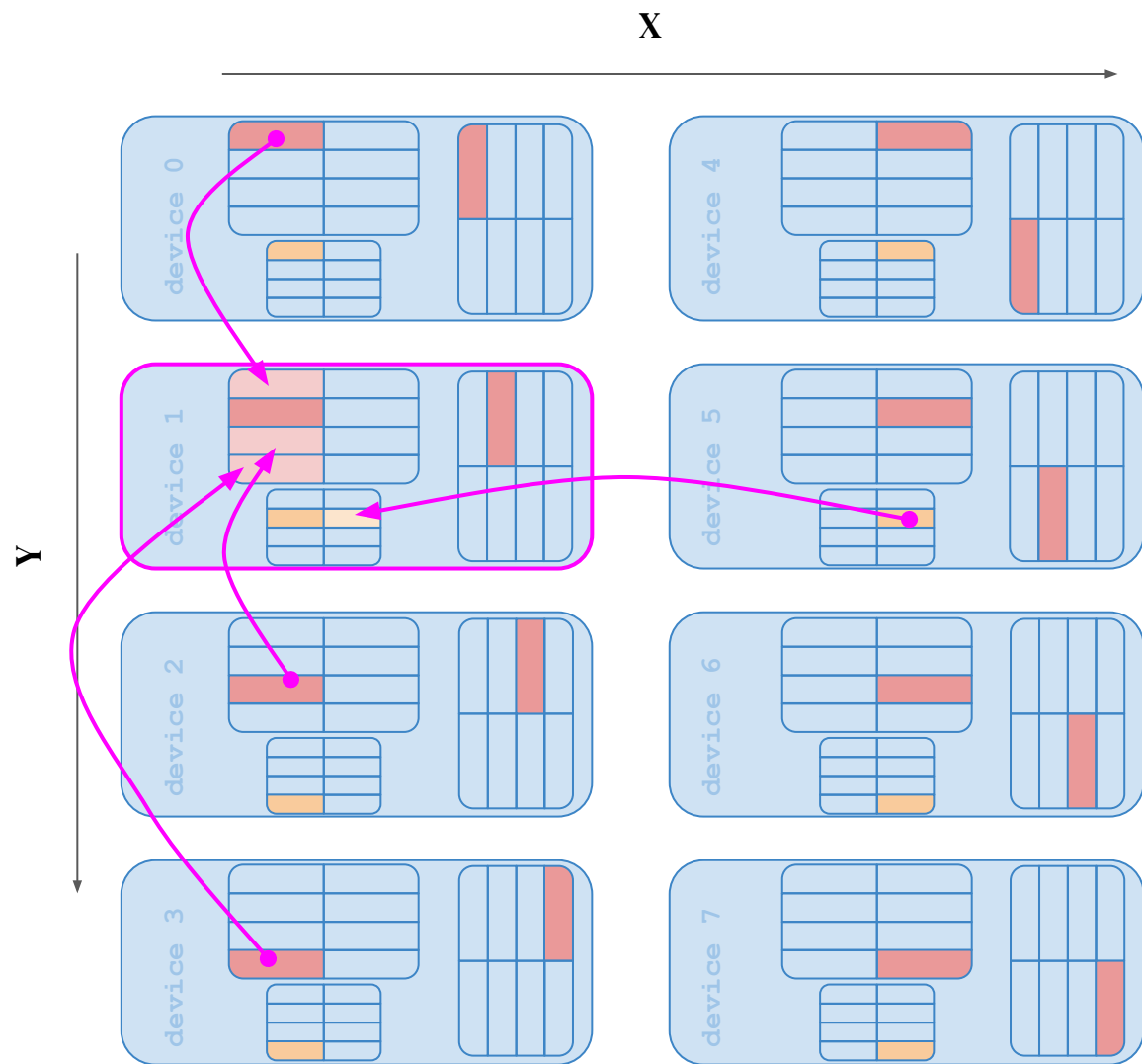
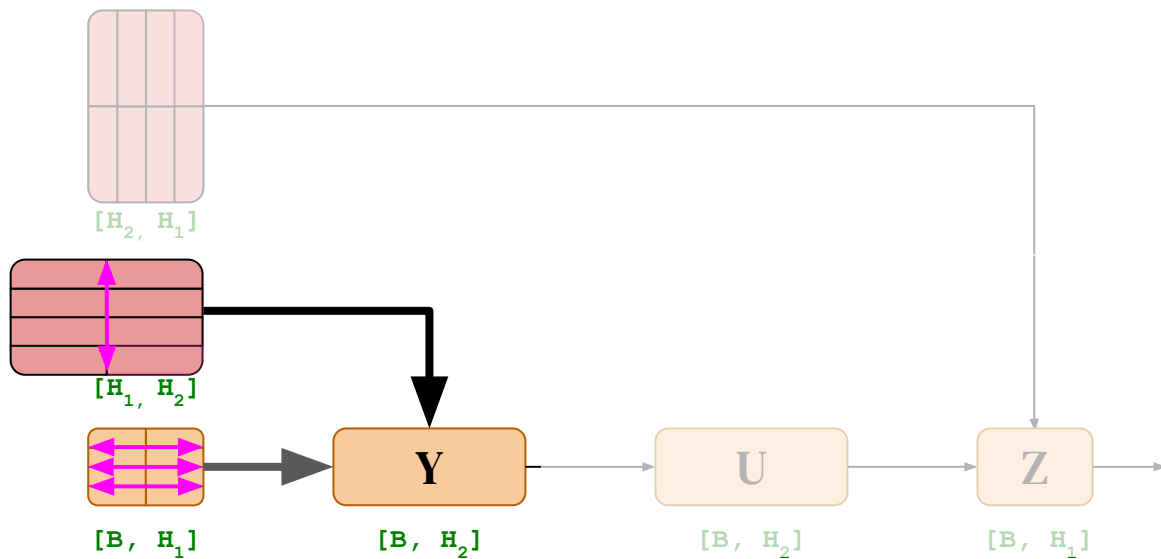
$[B/4, H_1/2] @ [H_1/4, H_2/2]$



AllGather \mathbf{D}_{yx} on X
 AllGather \mathbf{A}_{yx} on Y

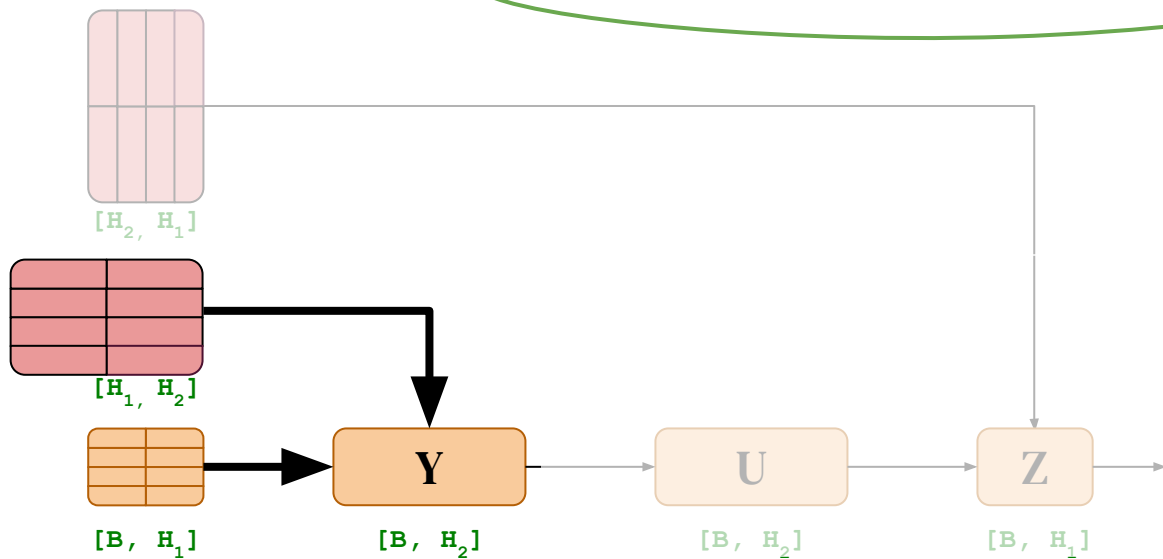


AllGather \mathbf{D}_{yx} on X
 AllGather \mathbf{A}_{yx} on Y



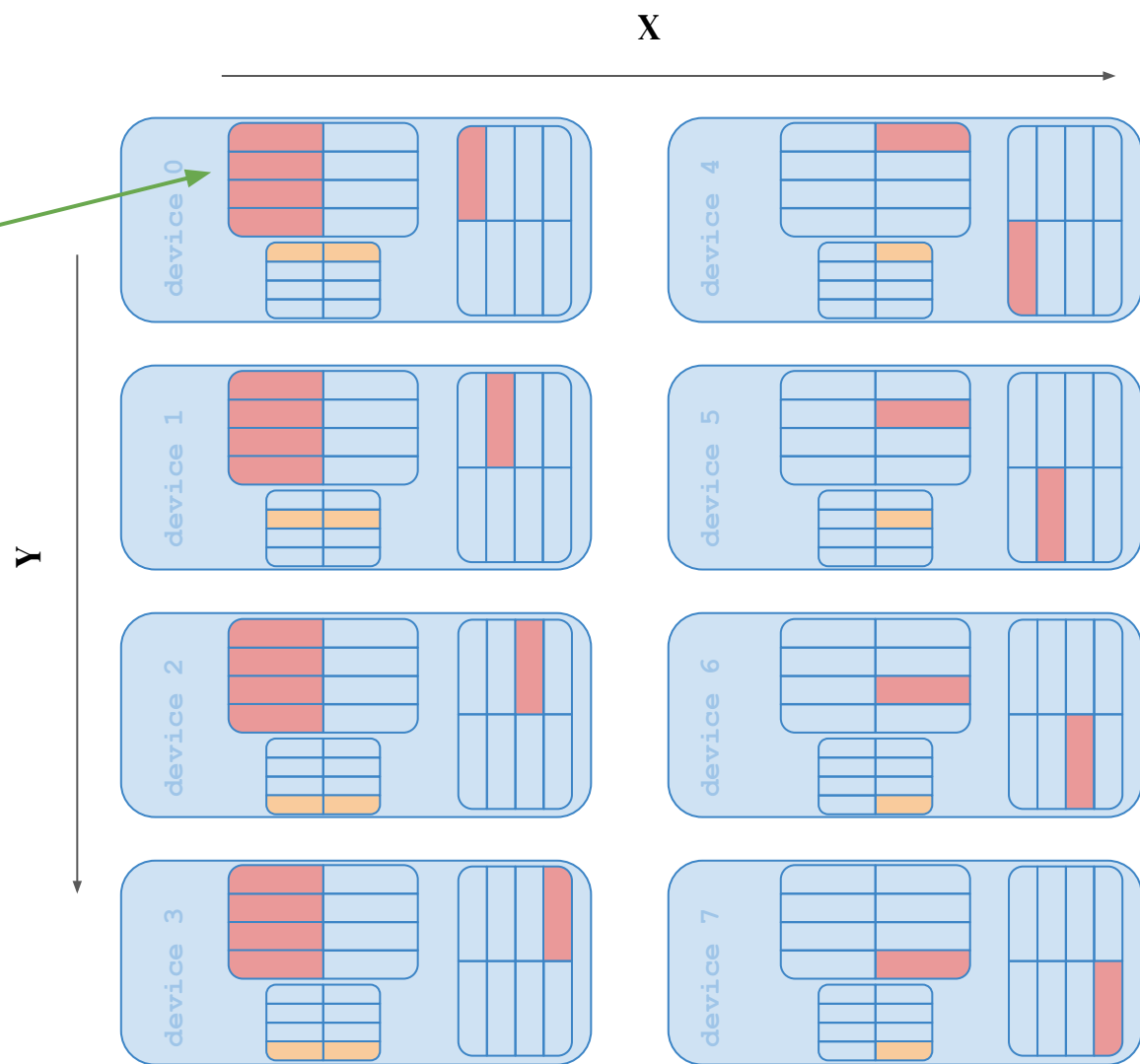
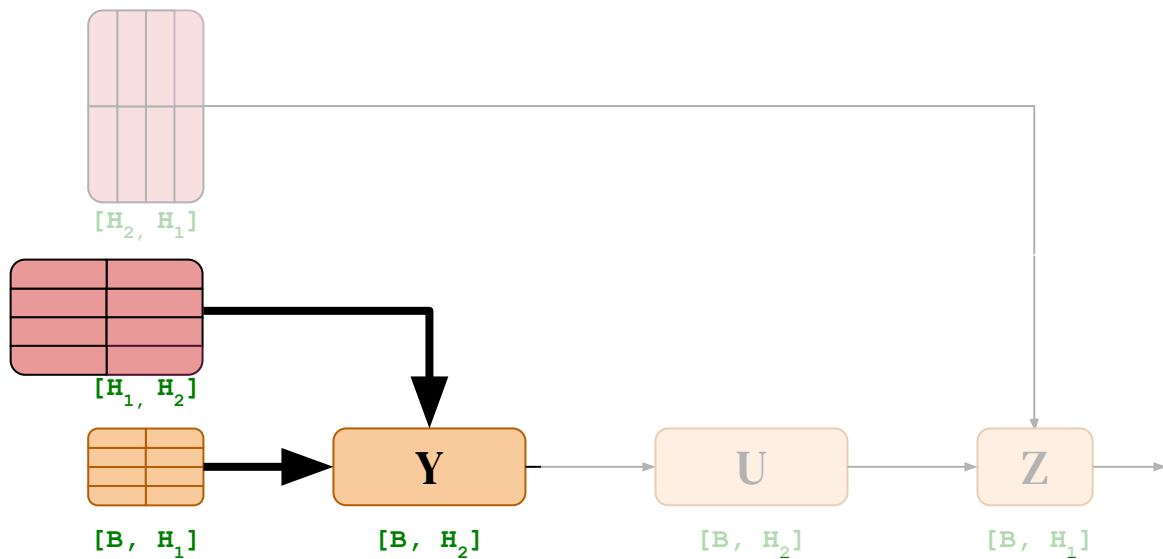
Compute $\mathbf{D}_{yX} @ \mathbf{A}_{YX}$

$[B/4, H_1]$



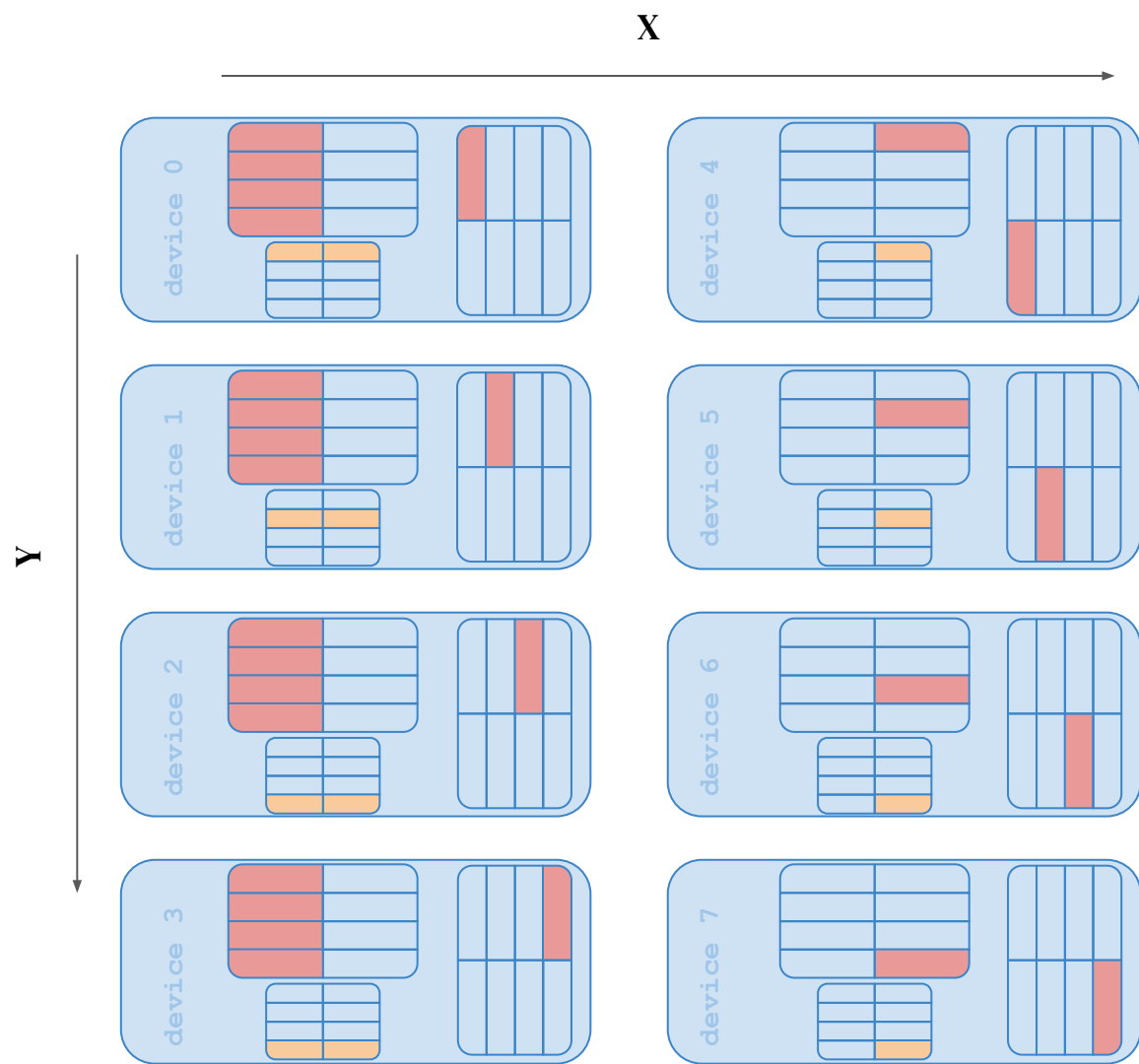
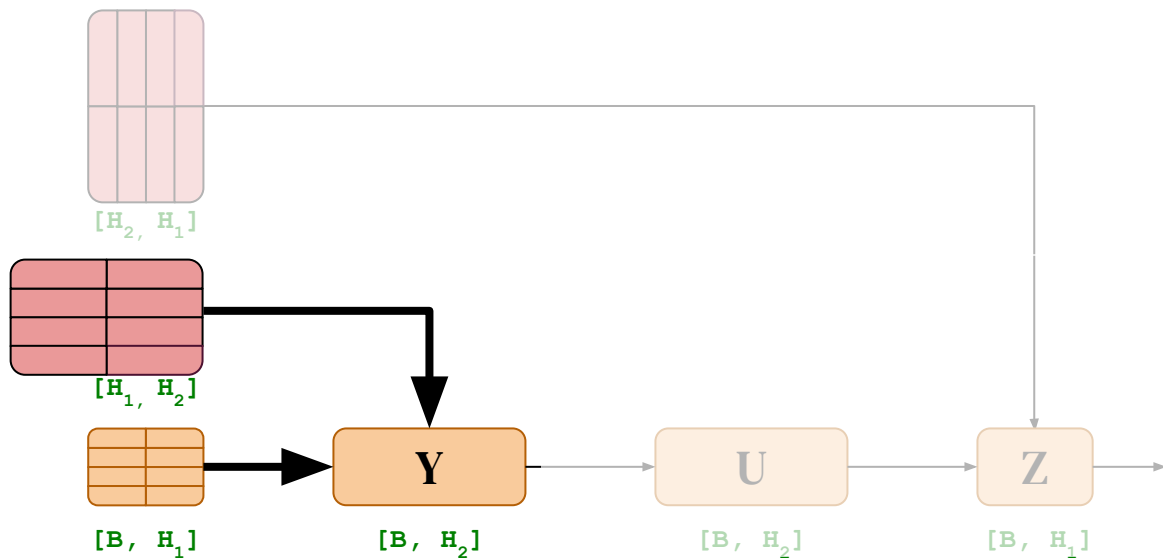
Compute $\mathbf{D}_{yX} @ \mathbf{A}_{YX}$

$[\mathbf{H}_1, \mathbf{H}_2/2]$



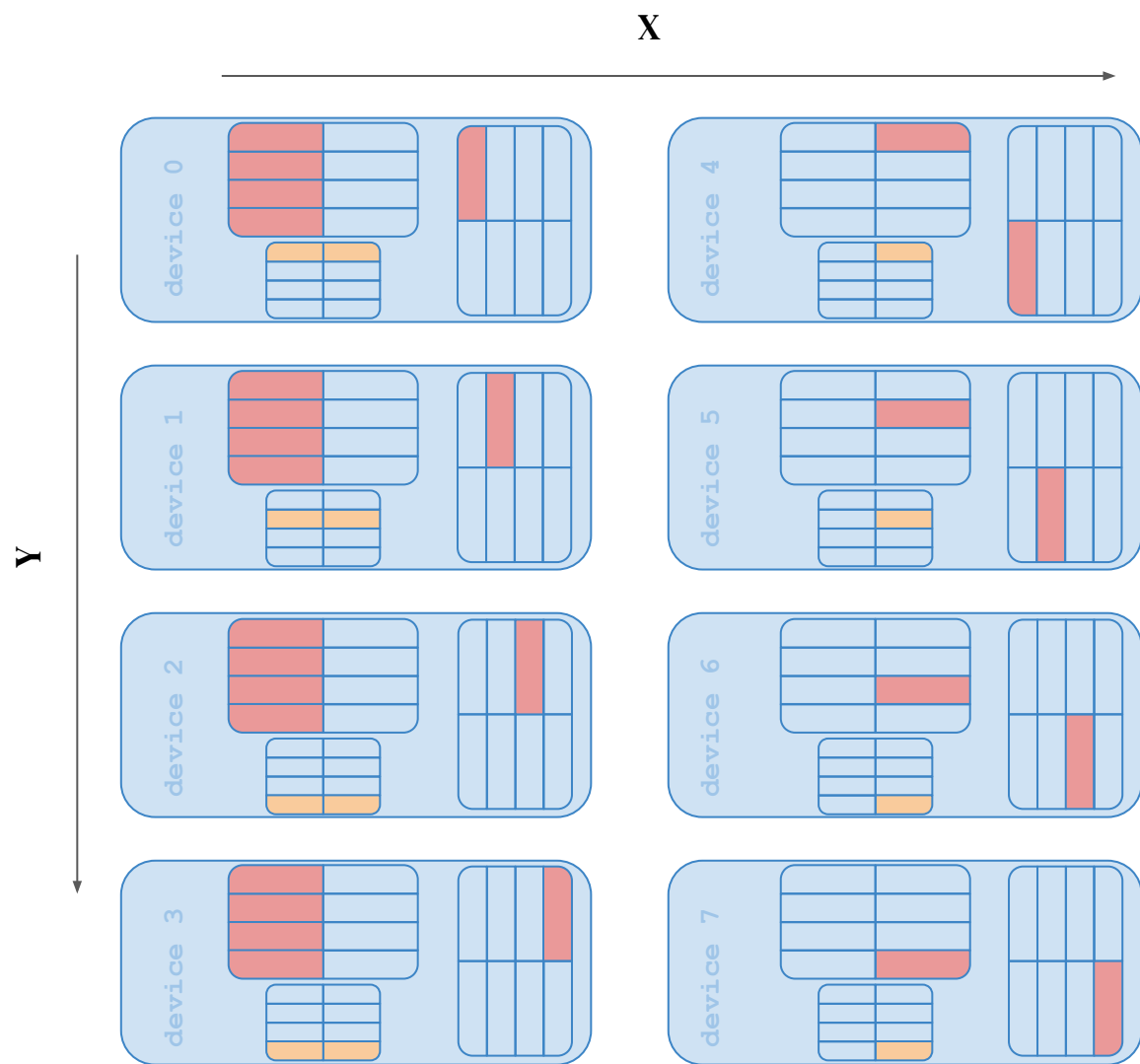
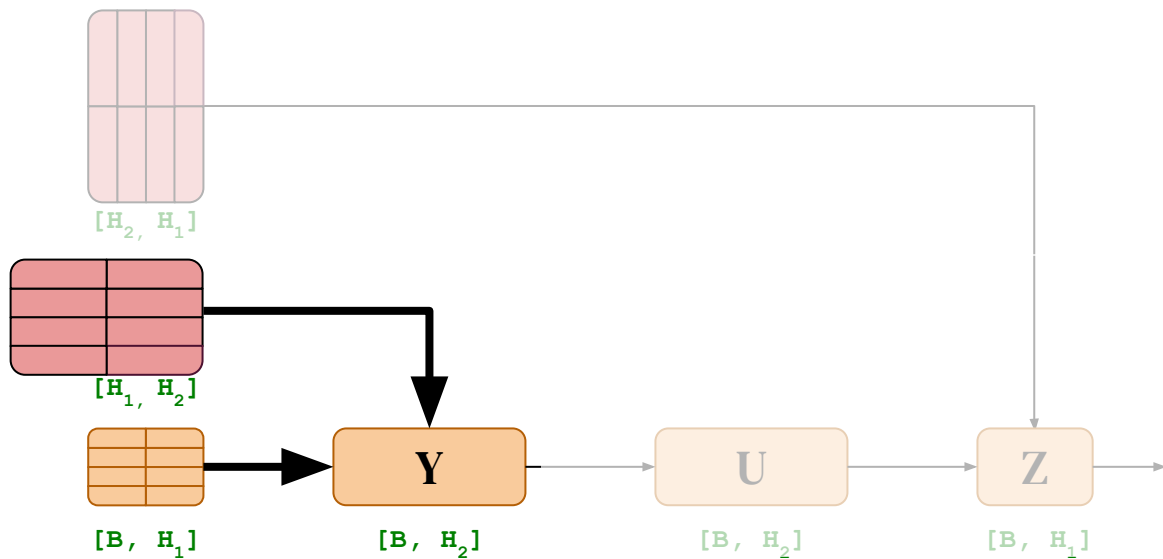
Compute $\mathbf{D}_{yX} @ \mathbf{A}_{YX}$

$$[B/4, H_1] @ [H_1, H_2/2]$$

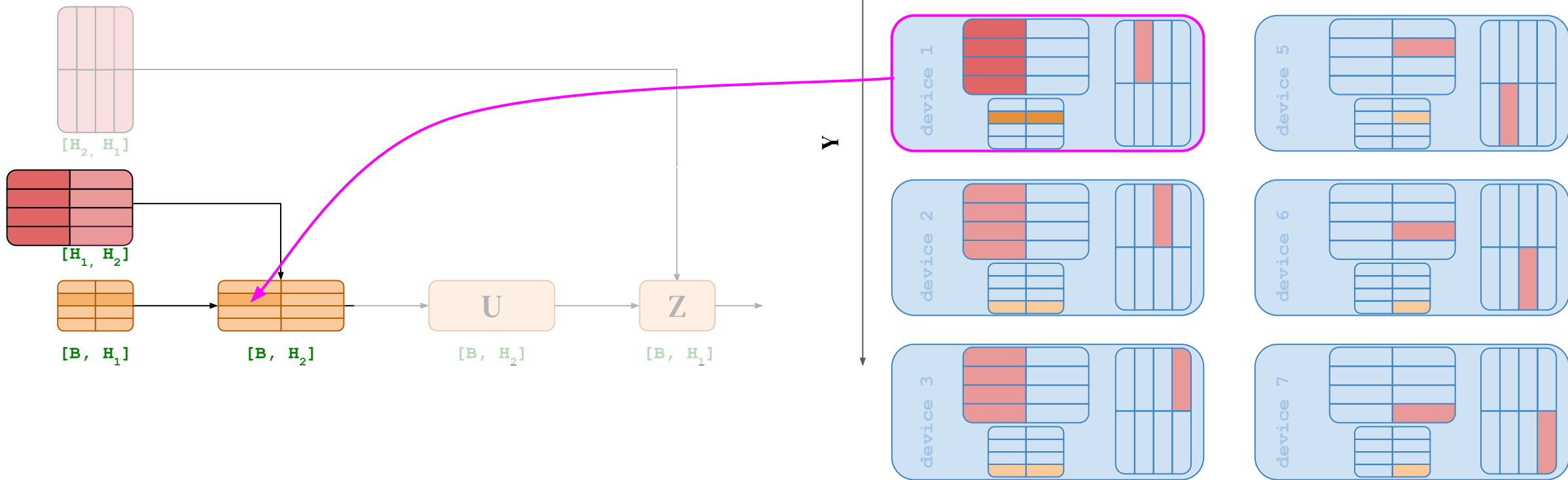


Compute $\mathbf{D}_{yX} @ \mathbf{A}_{YX}$

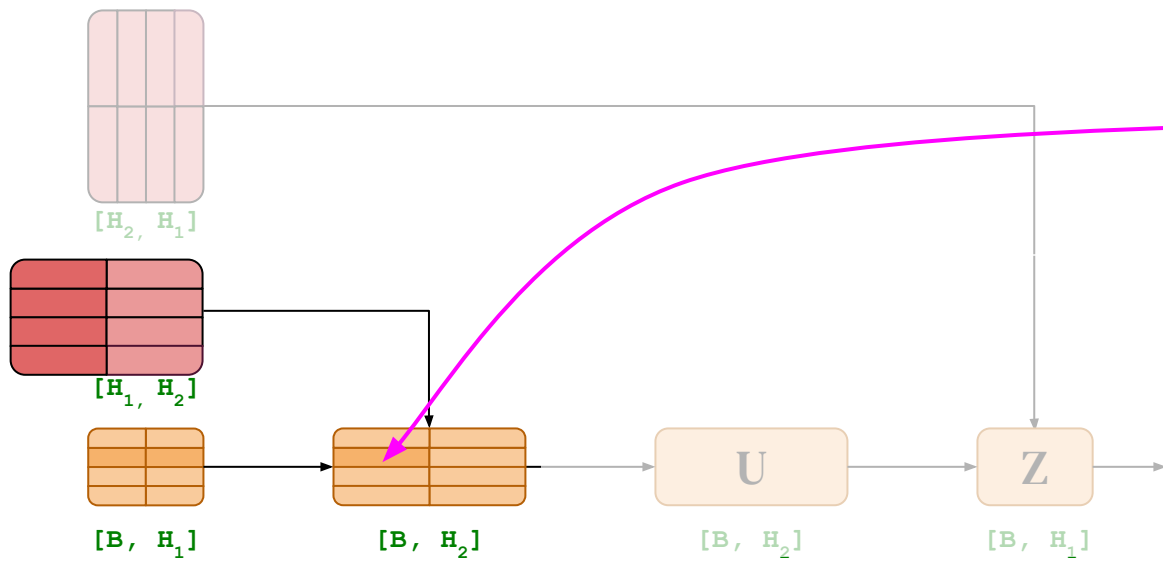
$$[B/4, H_1] @ [H_1, H_2/2]$$



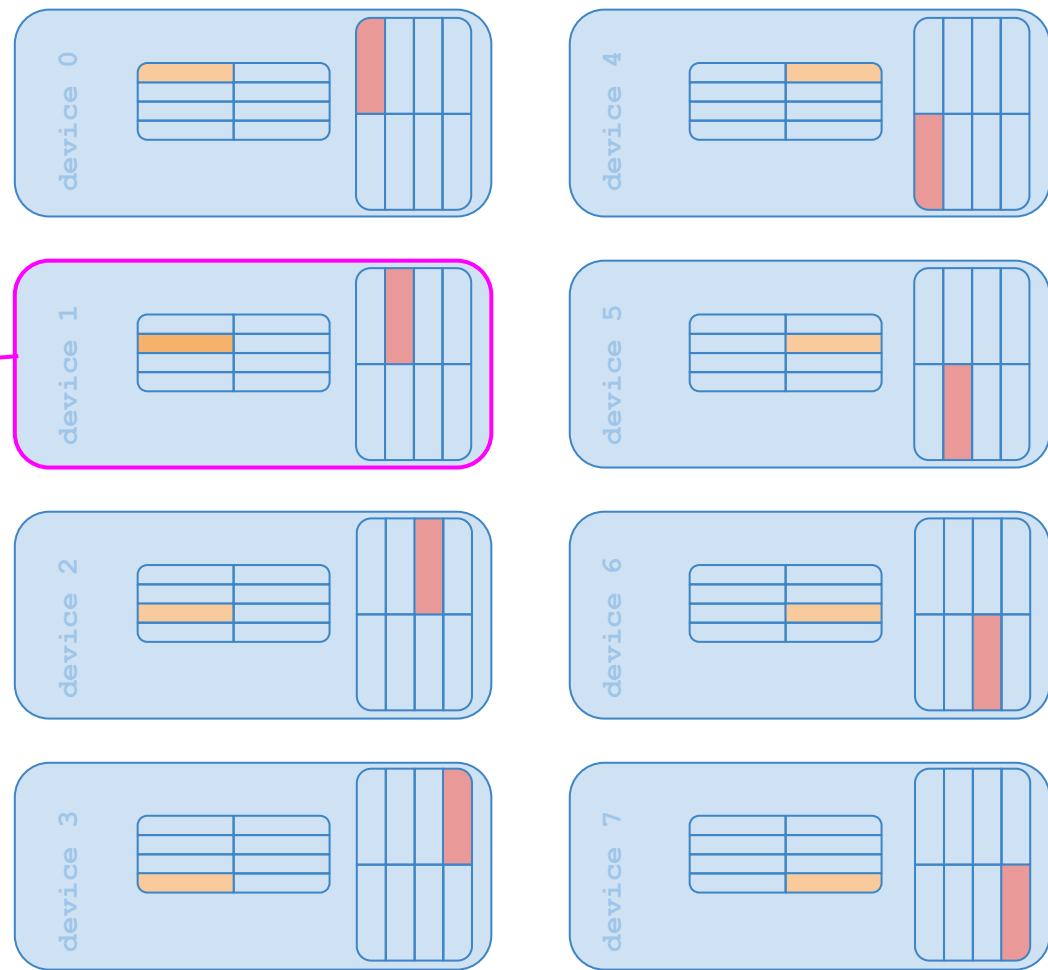
Compute $\mathbf{D}_{yX} @ \mathbf{A}_{YX}$



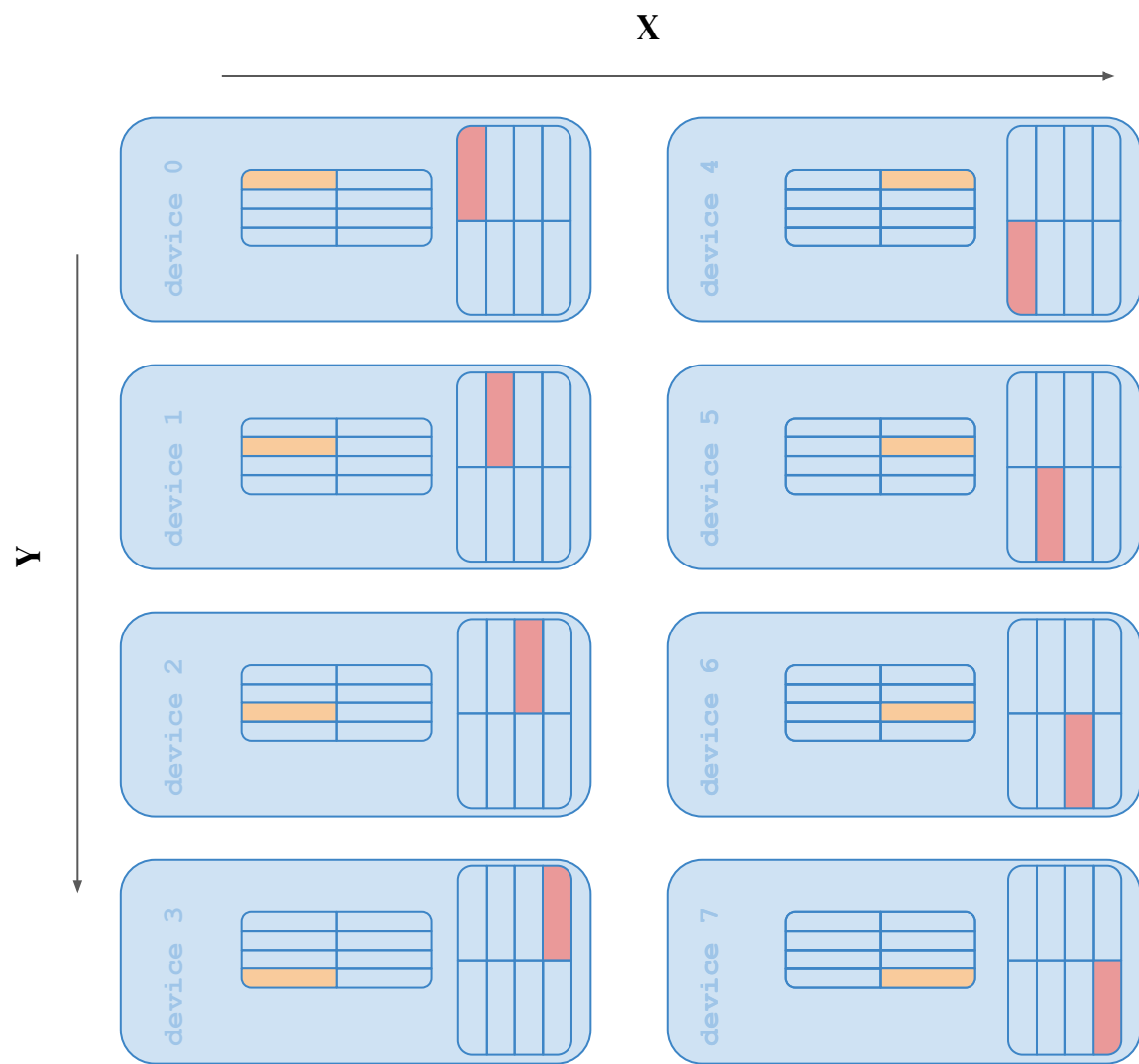
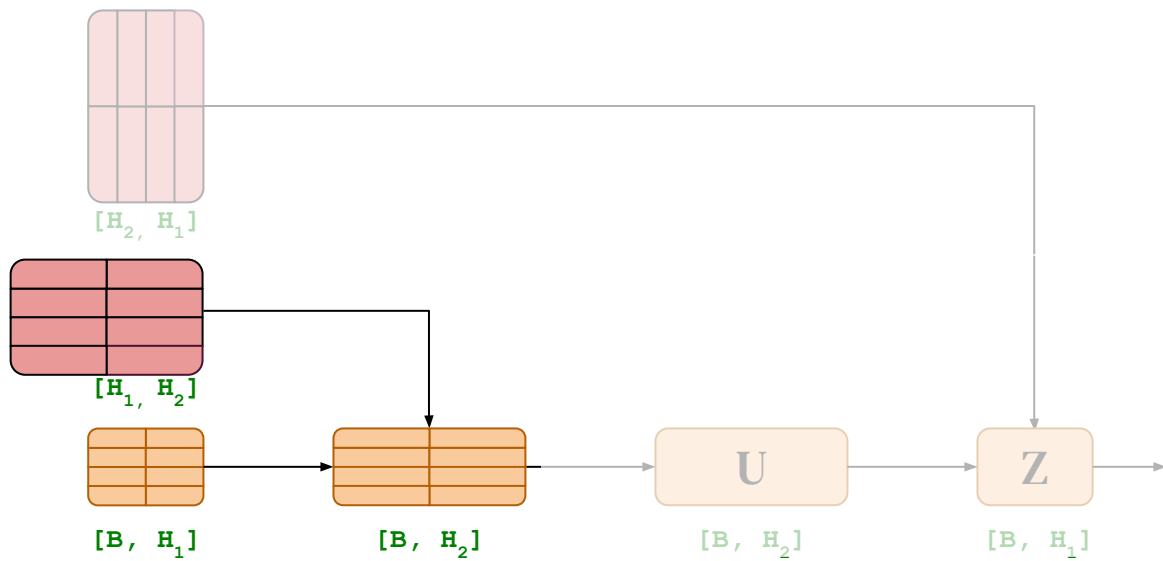
\mathbf{Y}_{yx}



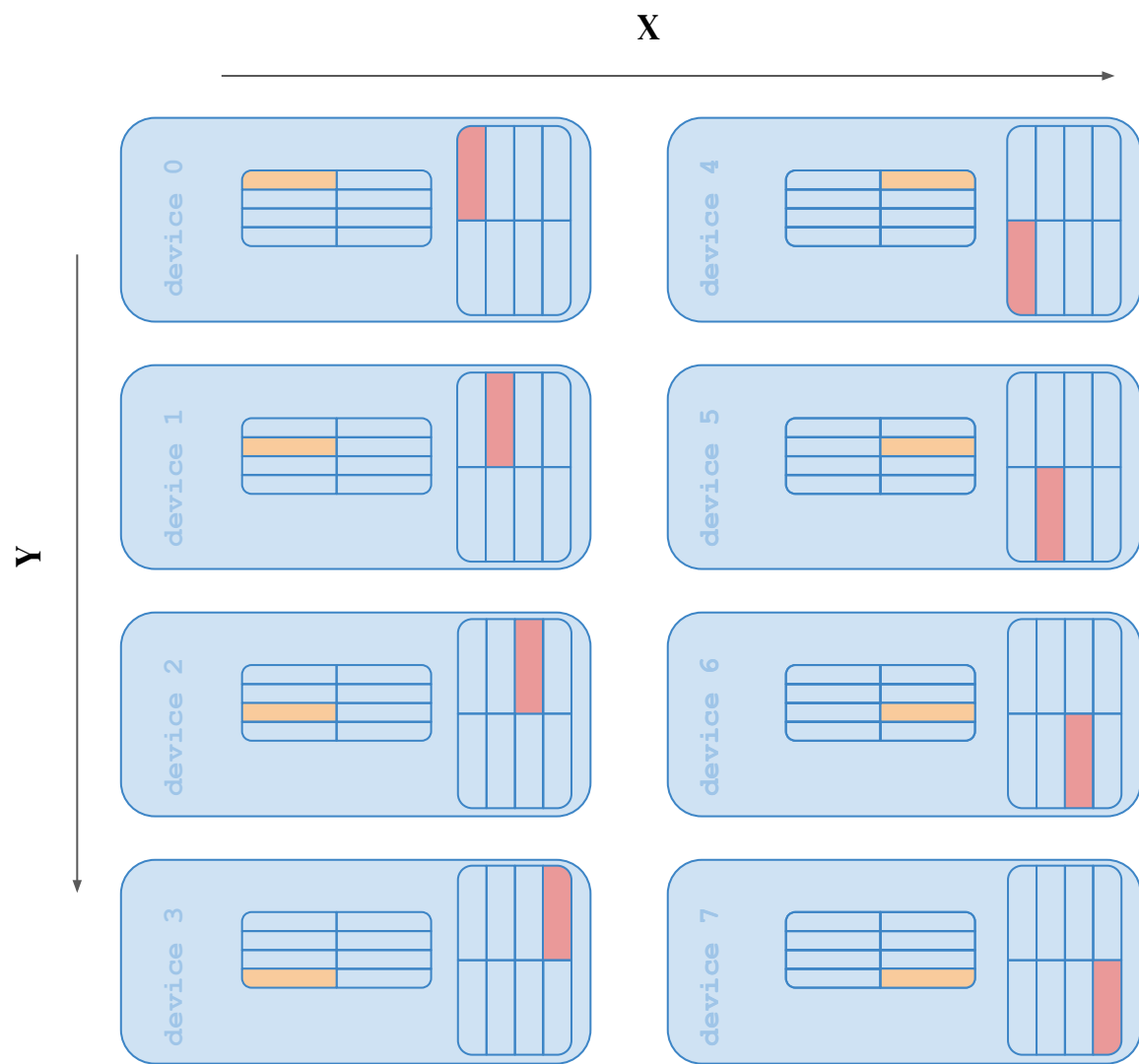
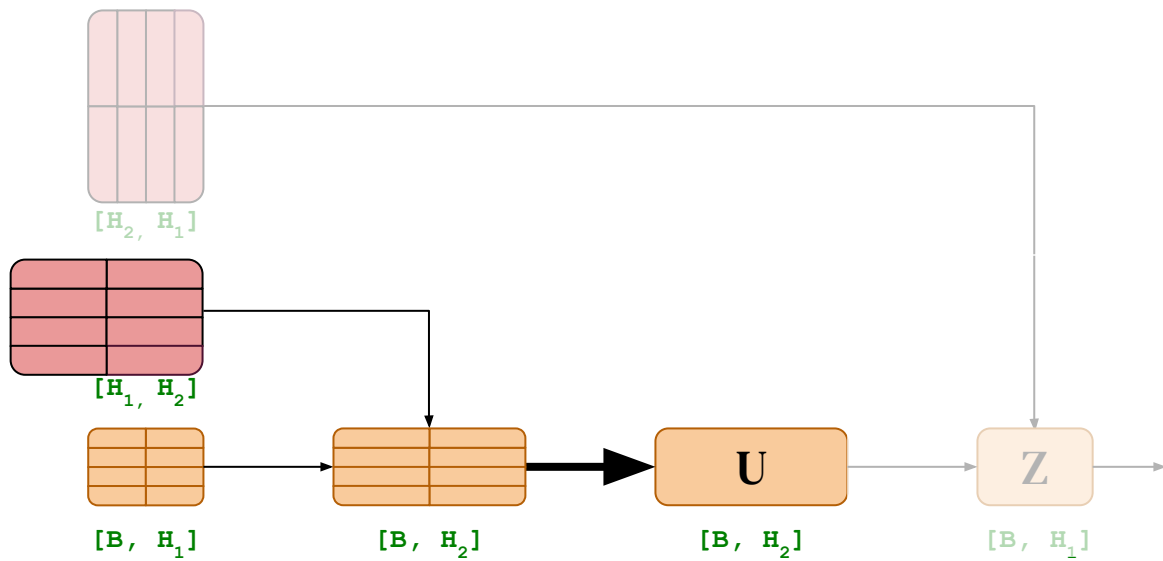
\mathbf{Y}



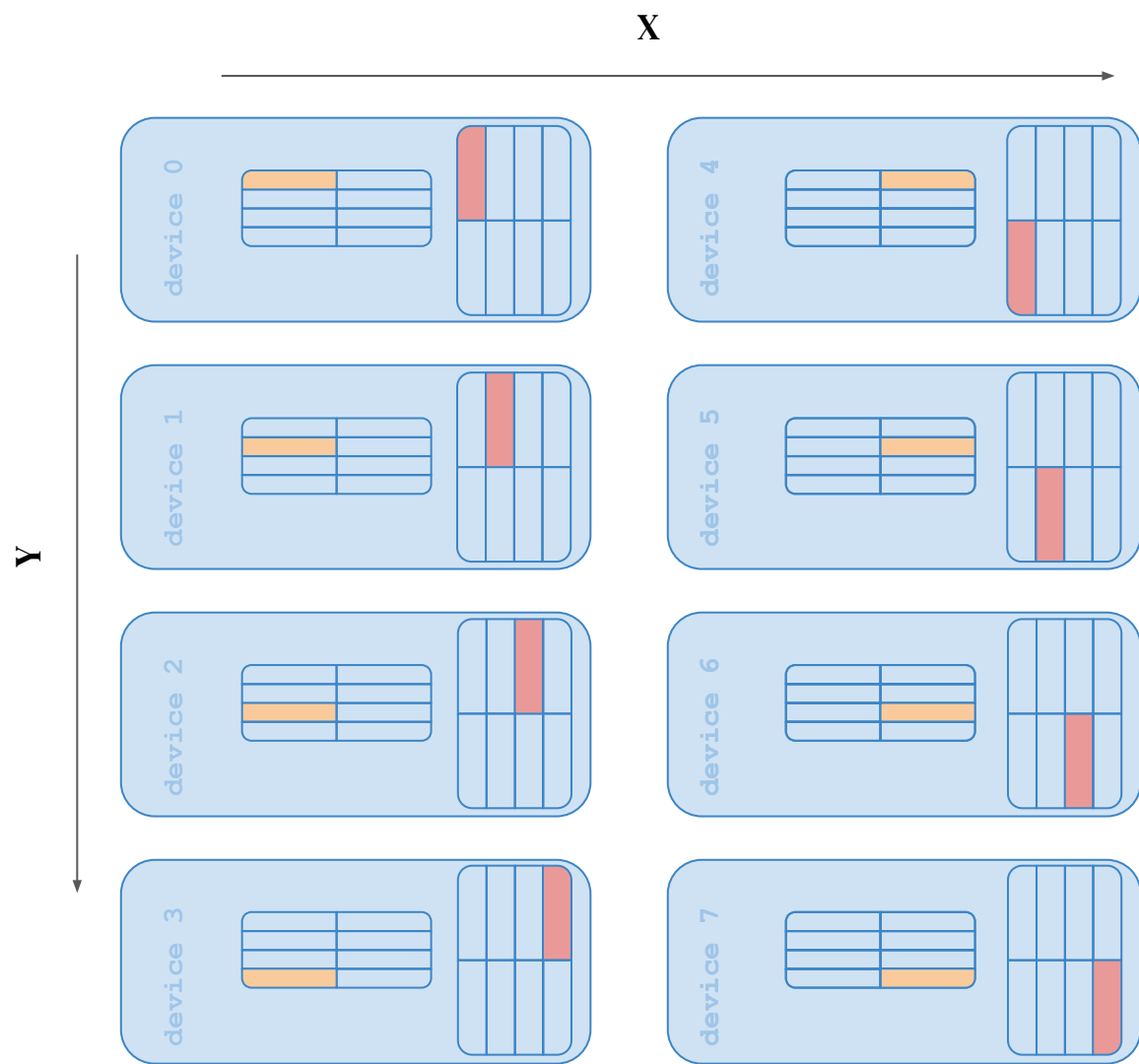
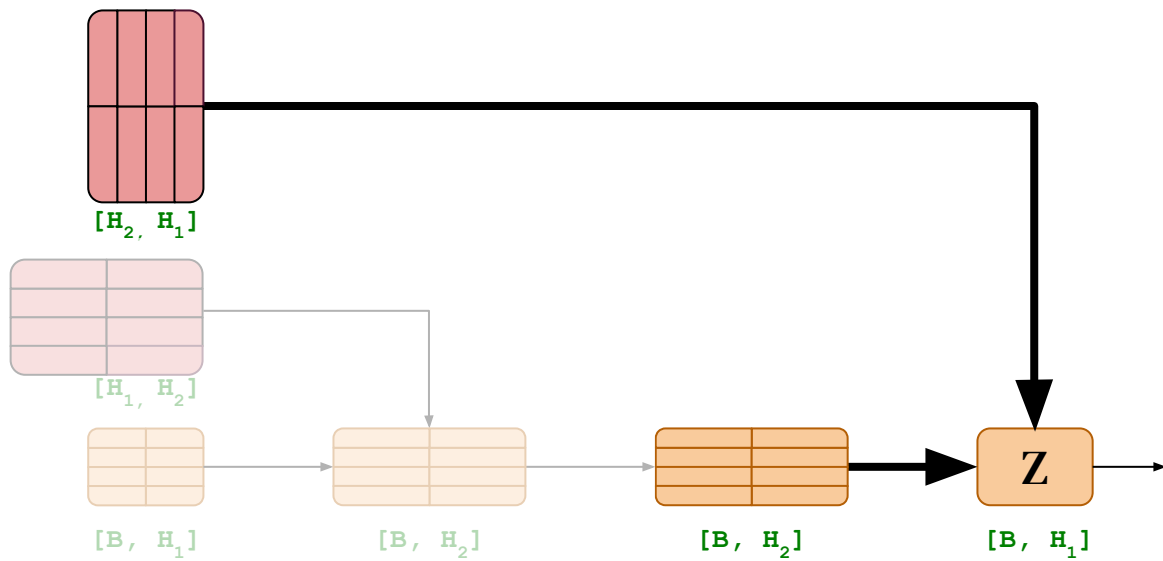
\mathbf{Y}_{yx}



Apply ReLU

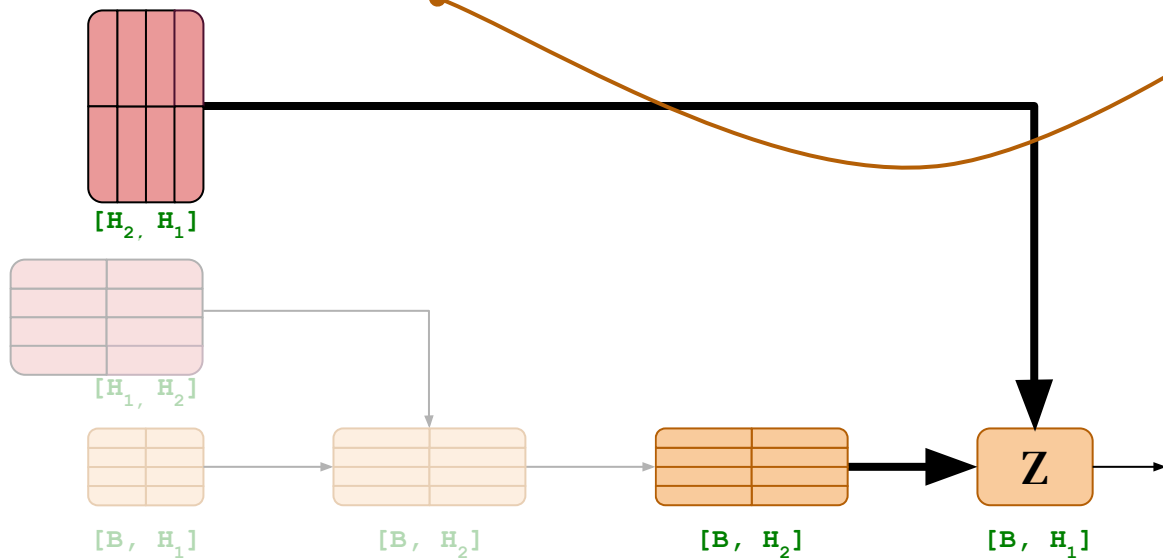


Compte \mathbf{U}_{yx} @ \mathbf{B}_{xy}



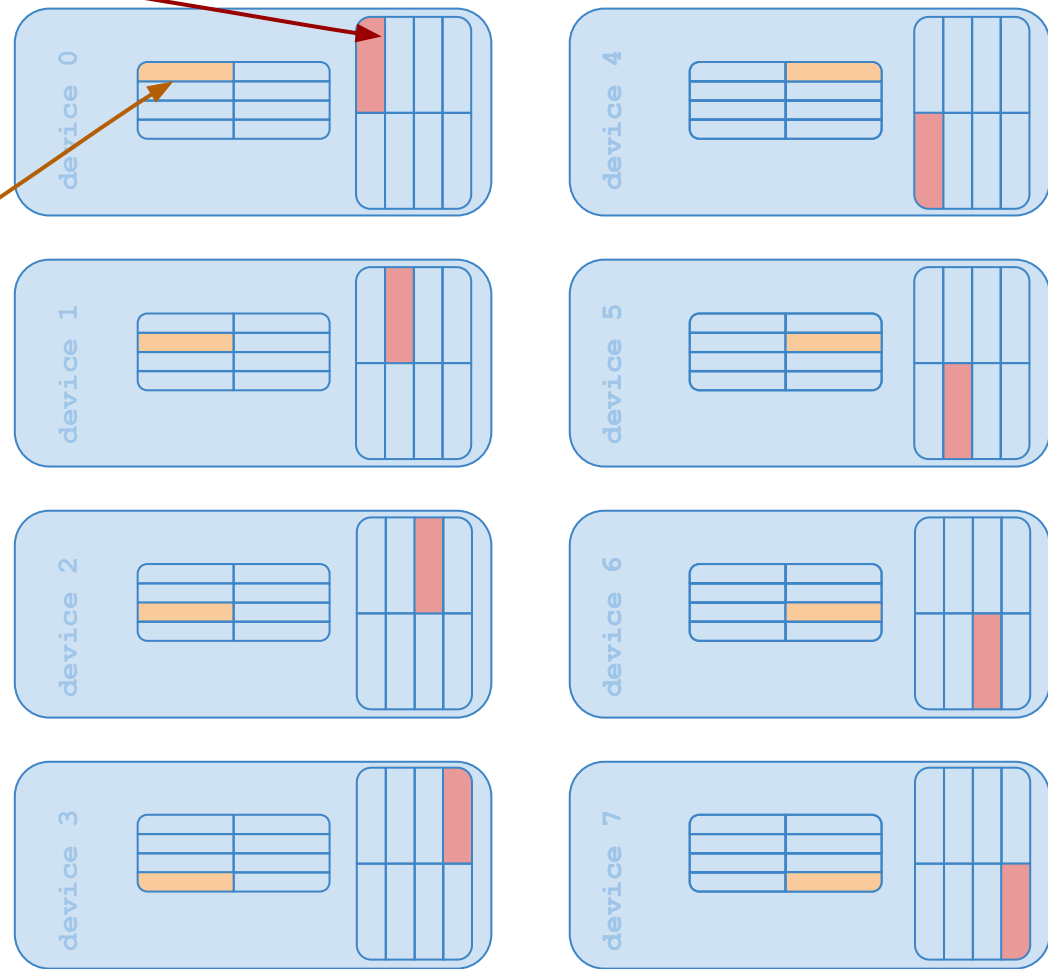
Compute $\mathbf{U}_{yx} @ \mathbf{B}_{xy}$

$[\mathbf{B}/4, \mathbf{H}_2/2] @ [\mathbf{H}_2/2, \mathbf{H}_1/4]$



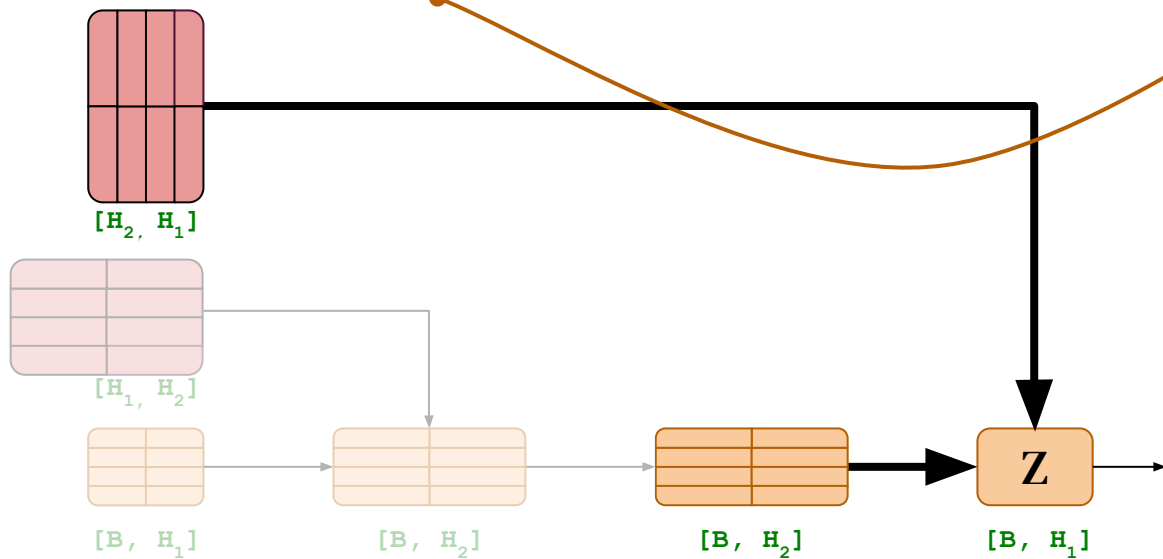
Y

X

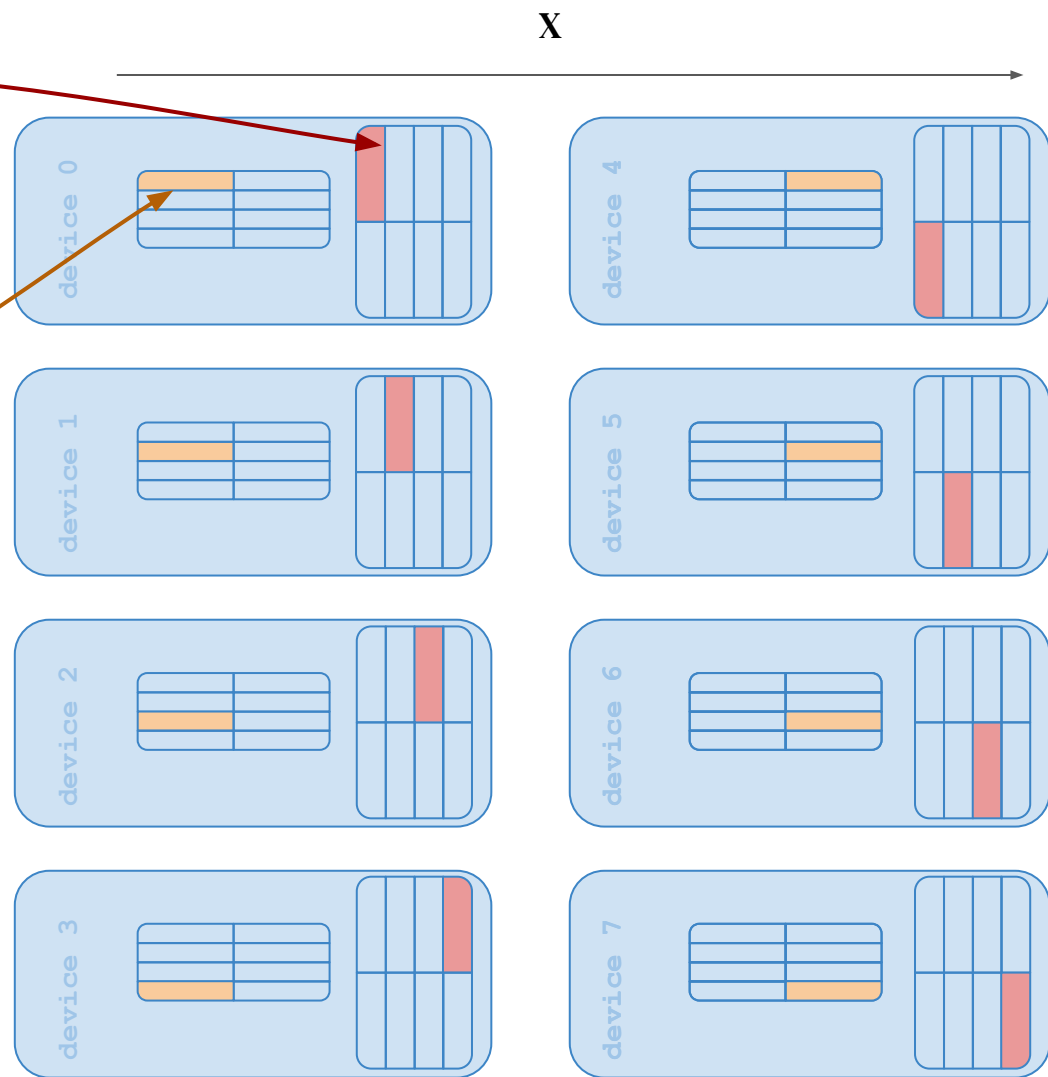


Compte \mathbf{U}_{yx} @ \mathbf{B}_{xy}

$[\mathbf{B}/4, \mathbf{H}_2/2] @ [\mathbf{H}_2/2, \mathbf{H}_1/4]$

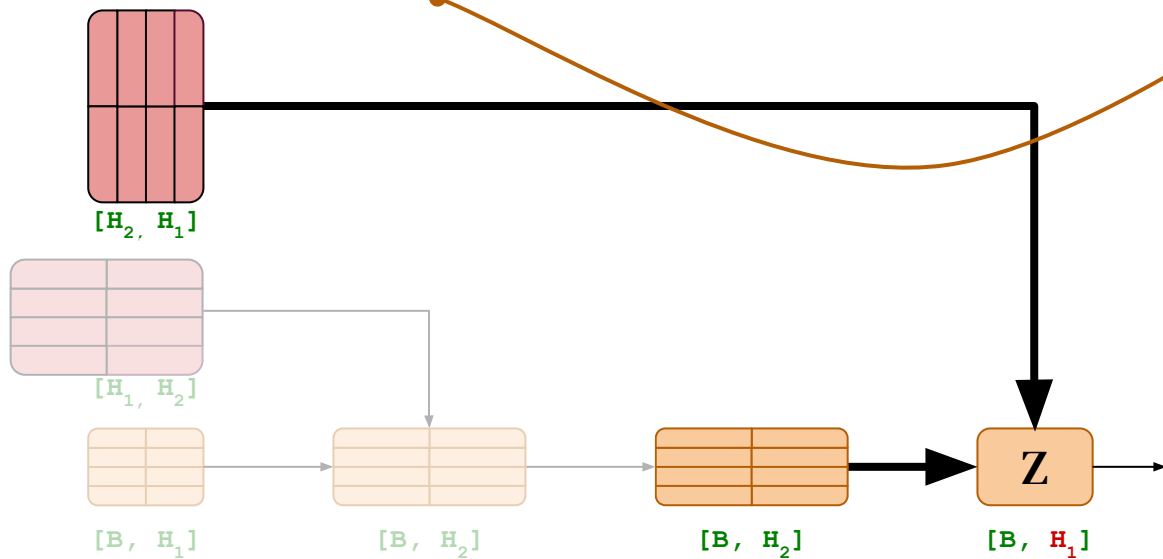


\mathbf{Y}



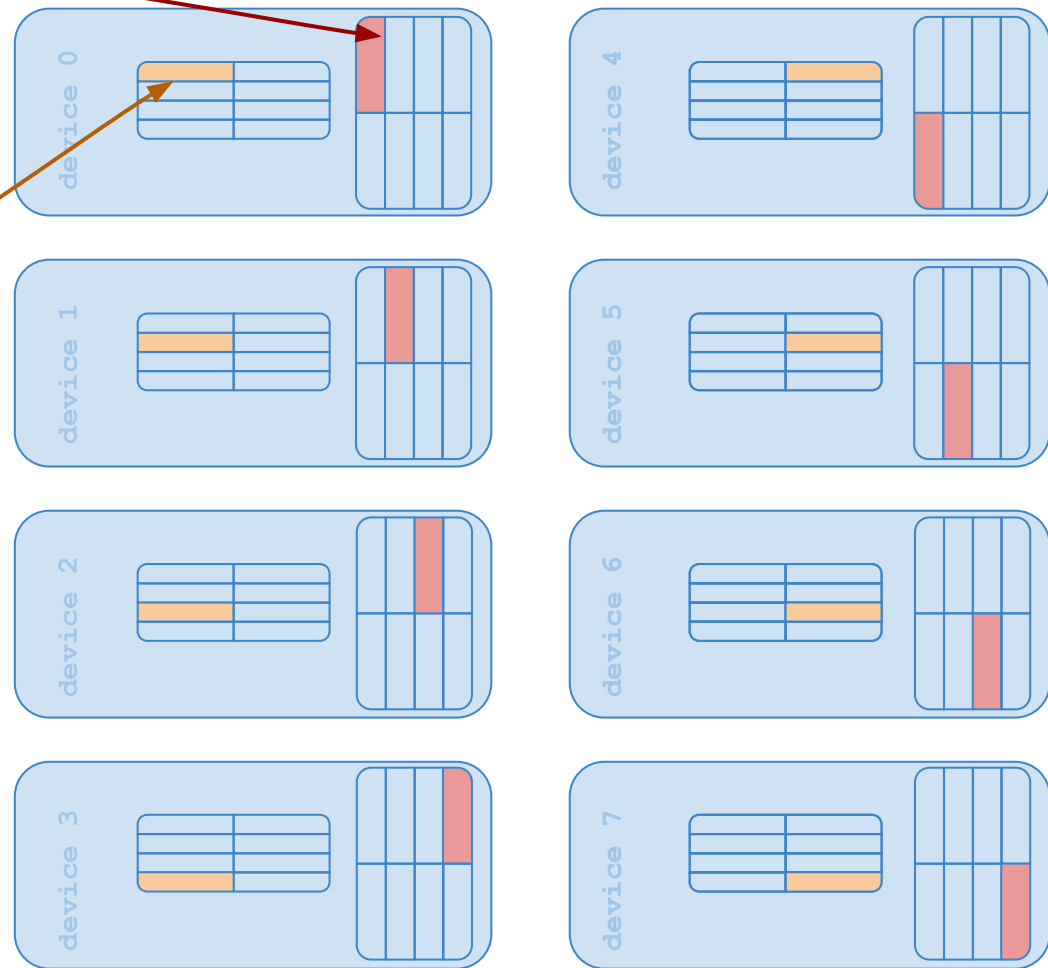
Compte \mathbf{U}_{yx} @ \mathbf{B}_{xy}

$[\mathbf{B}/4, \mathbf{H}_2/2] @ [\mathbf{H}_2/2, \mathbf{H}_1/4]$

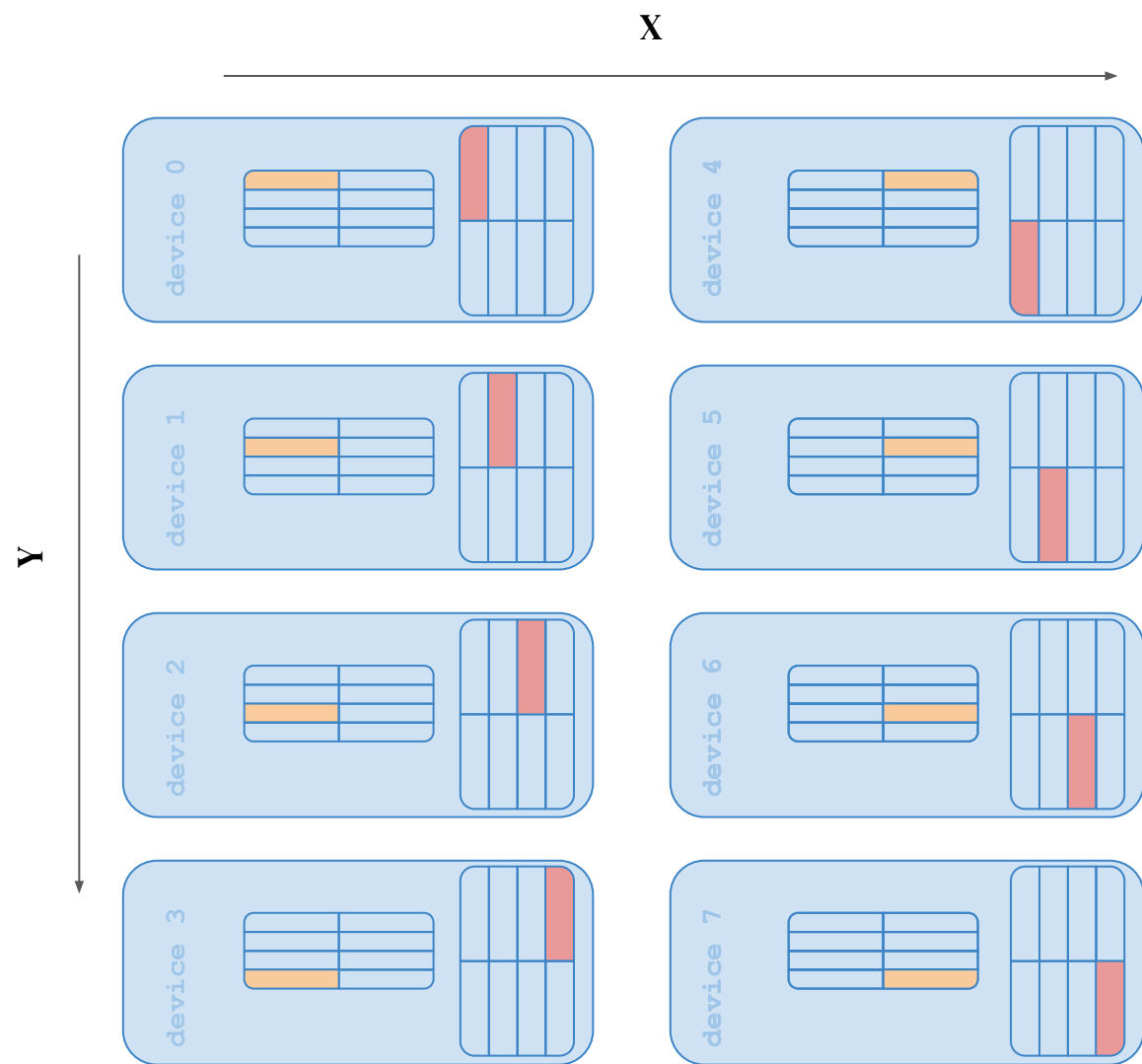
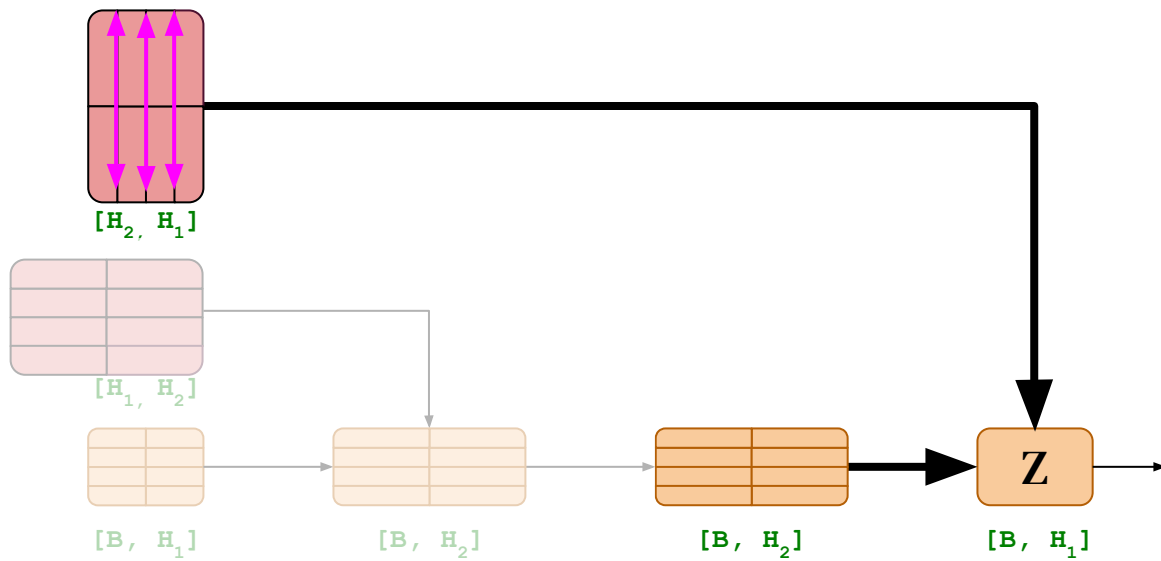


Y

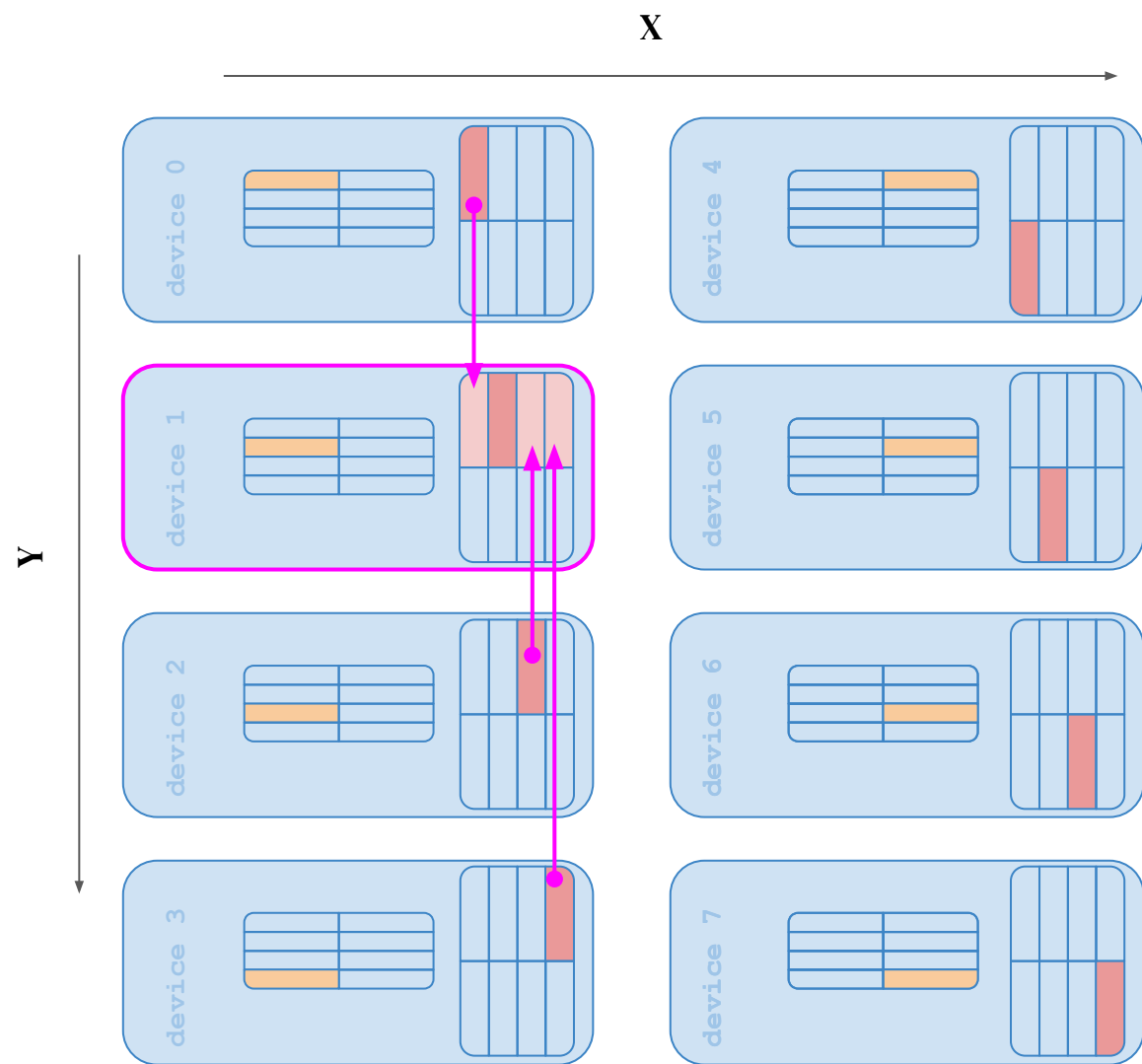
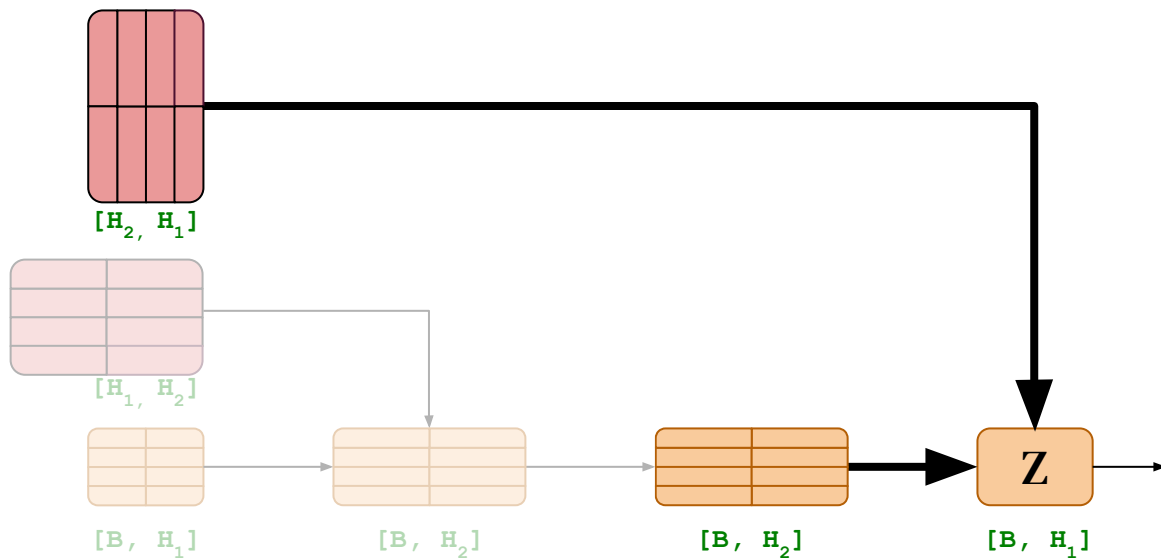
X



AllGather \mathbf{B}_{xy} on \mathbf{Y}

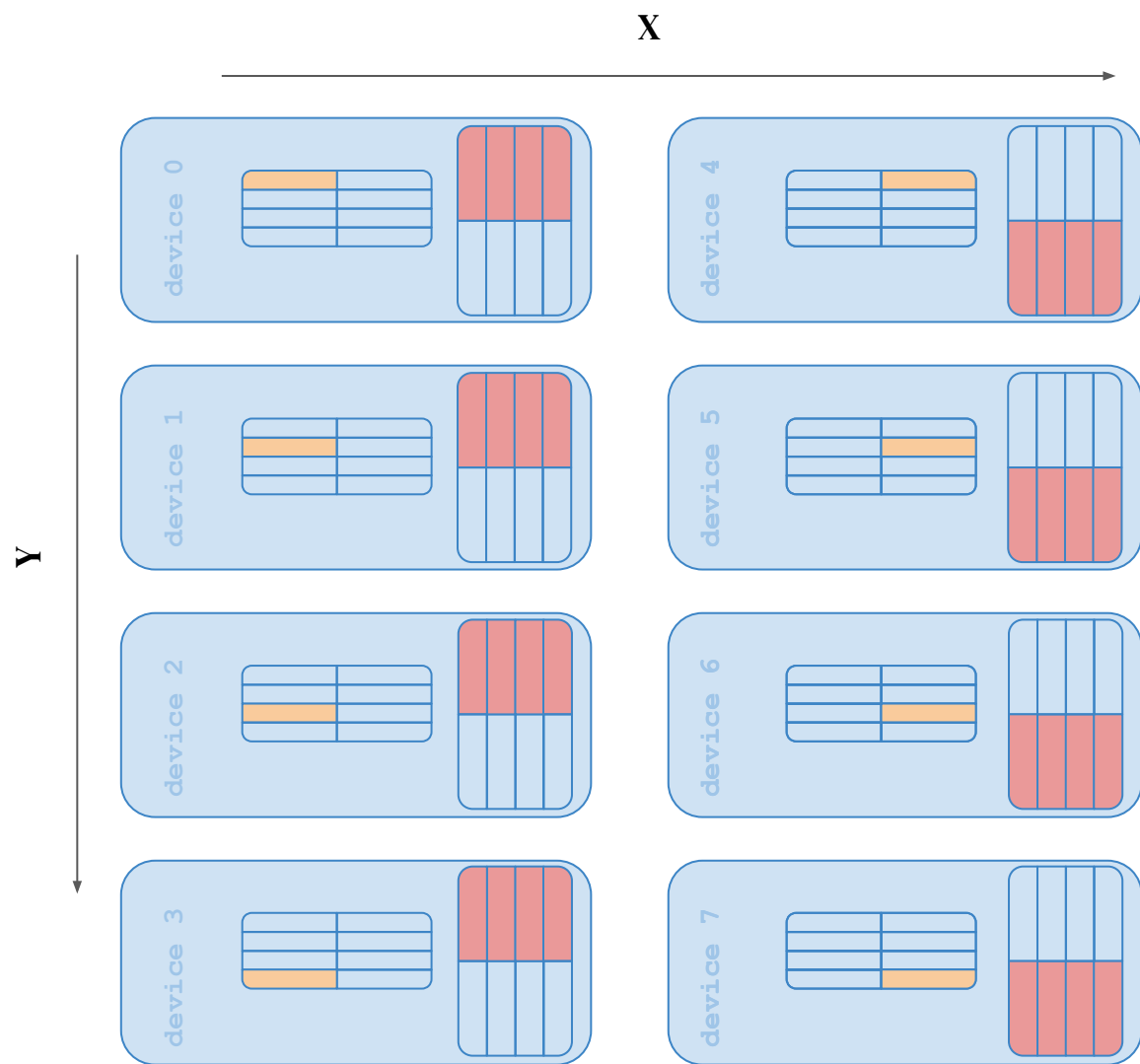
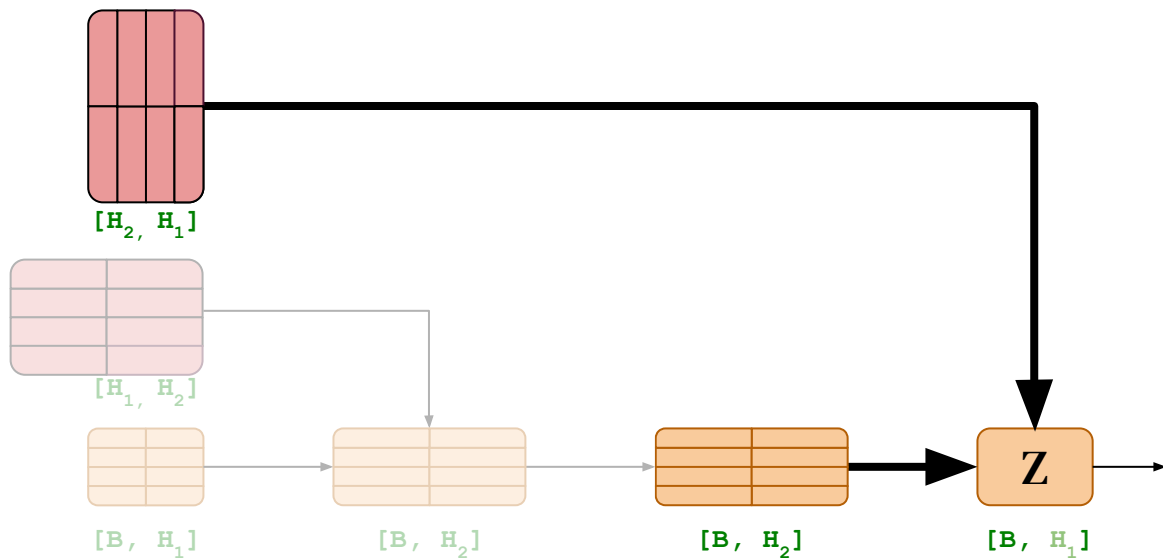


AllGather \mathbf{B}_{xy} on \mathbf{Y}

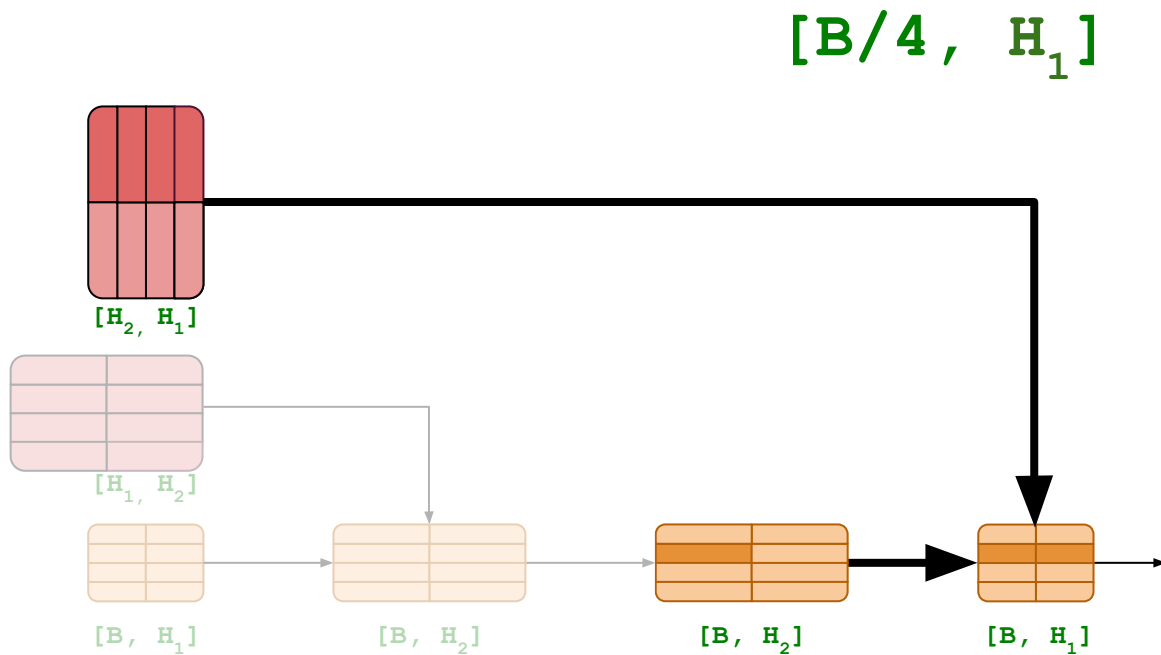


Compute $\mathbf{U}_{yx} @ \mathbf{B}_{xy}$

$$[\mathbf{B}/4, H_2/2] @ [H_2/2, H_1]$$

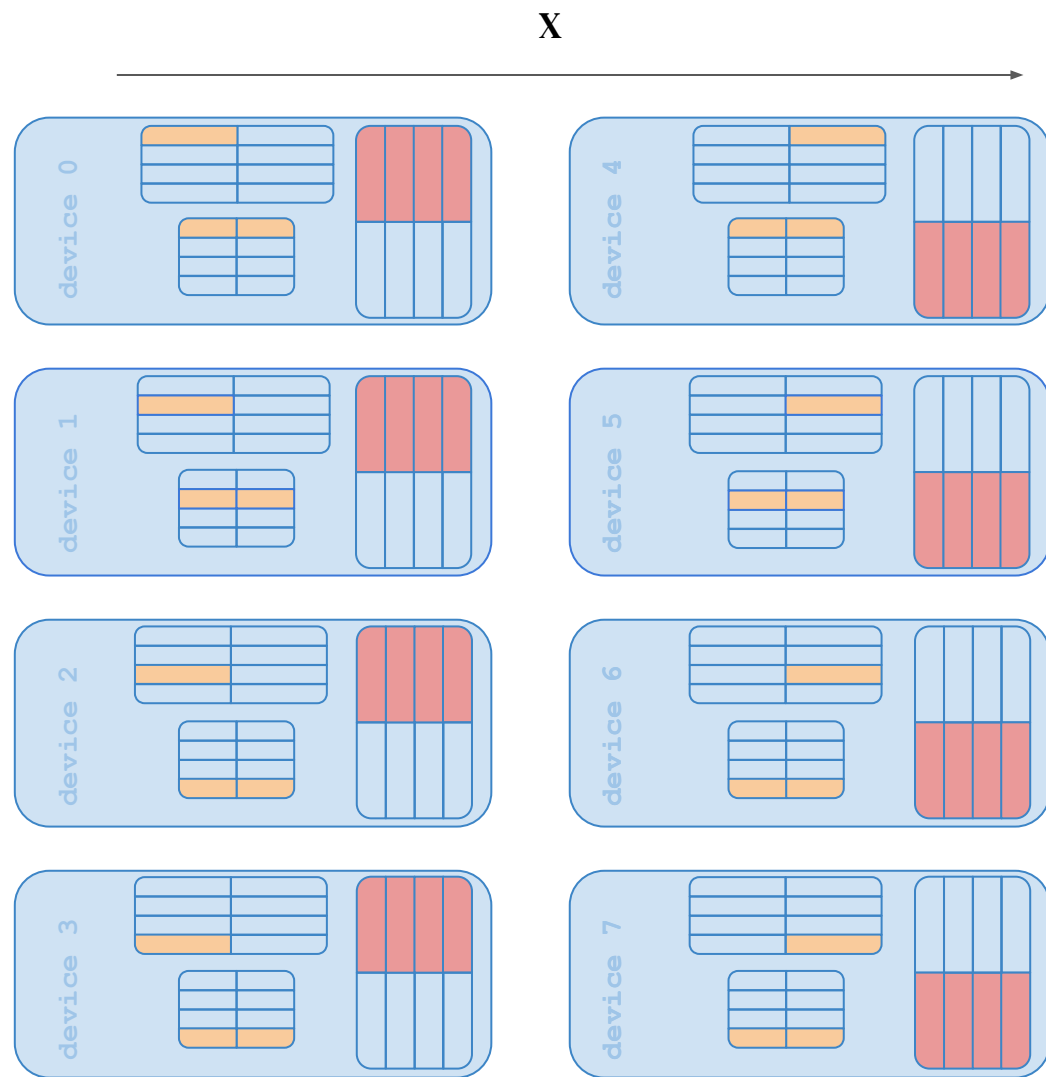


Are we done?

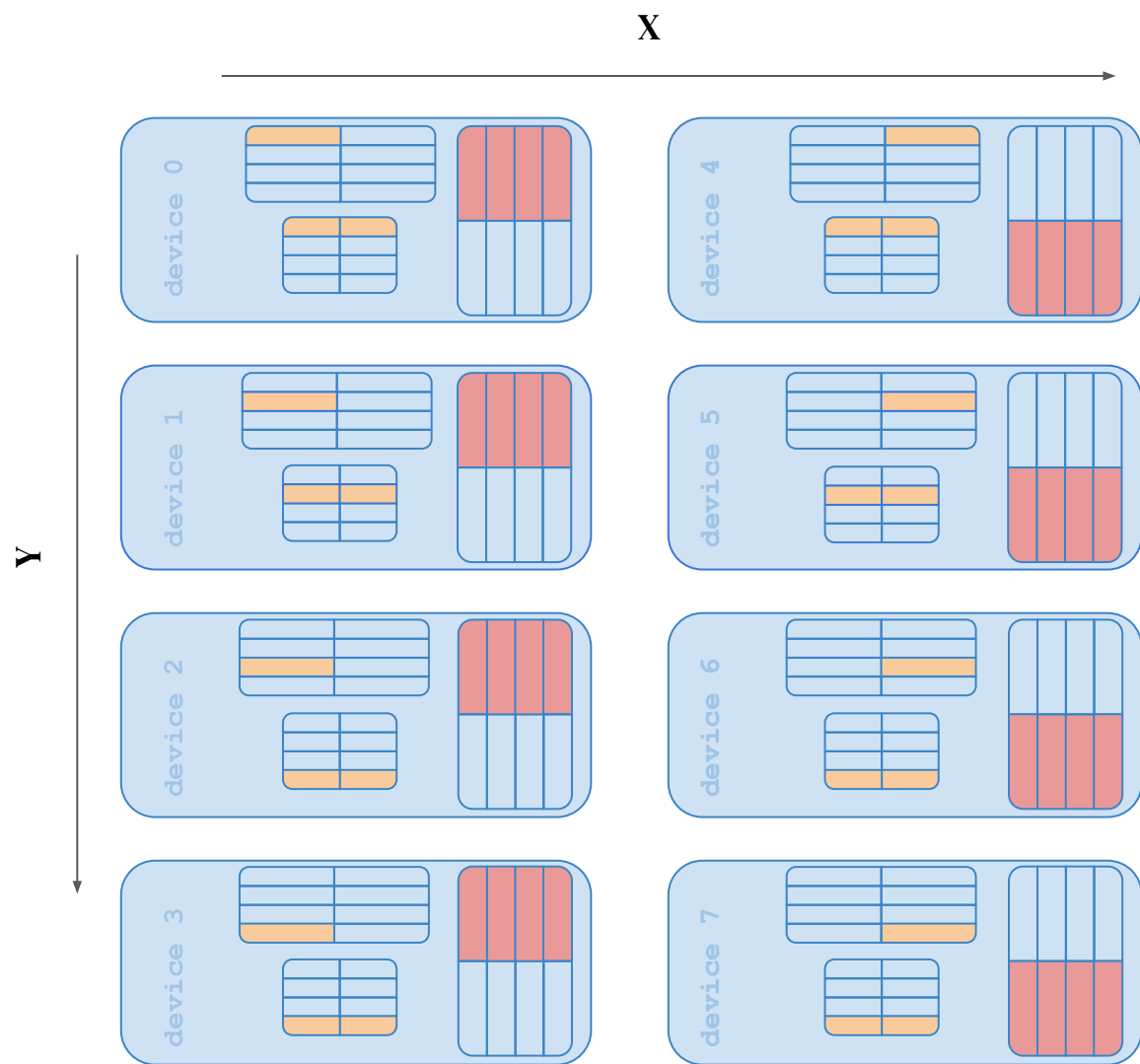
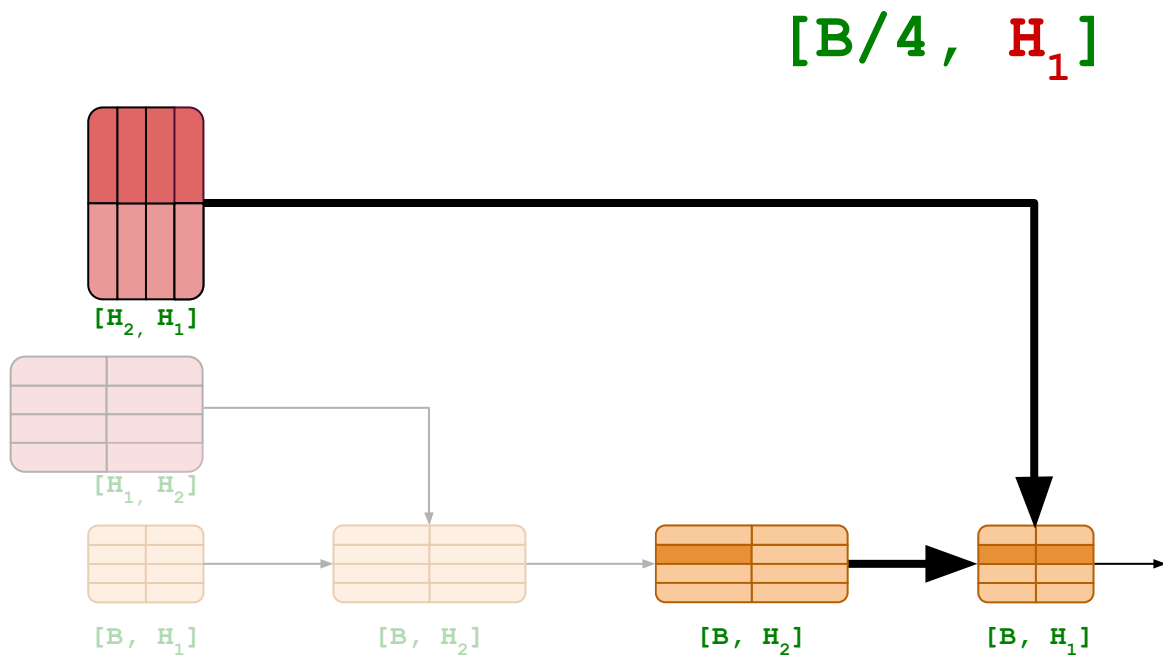


$[B/4, H_1]$

Y

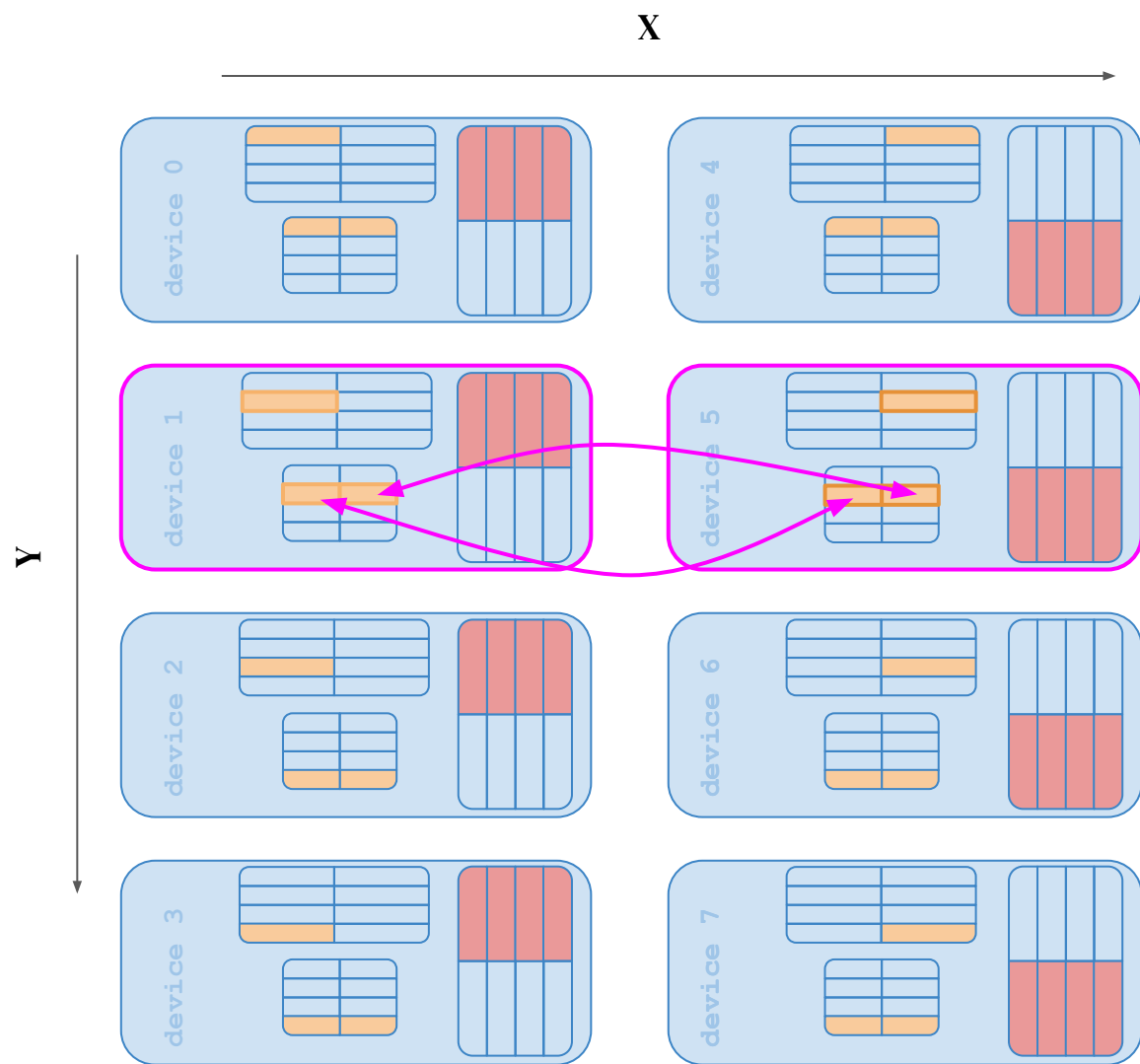
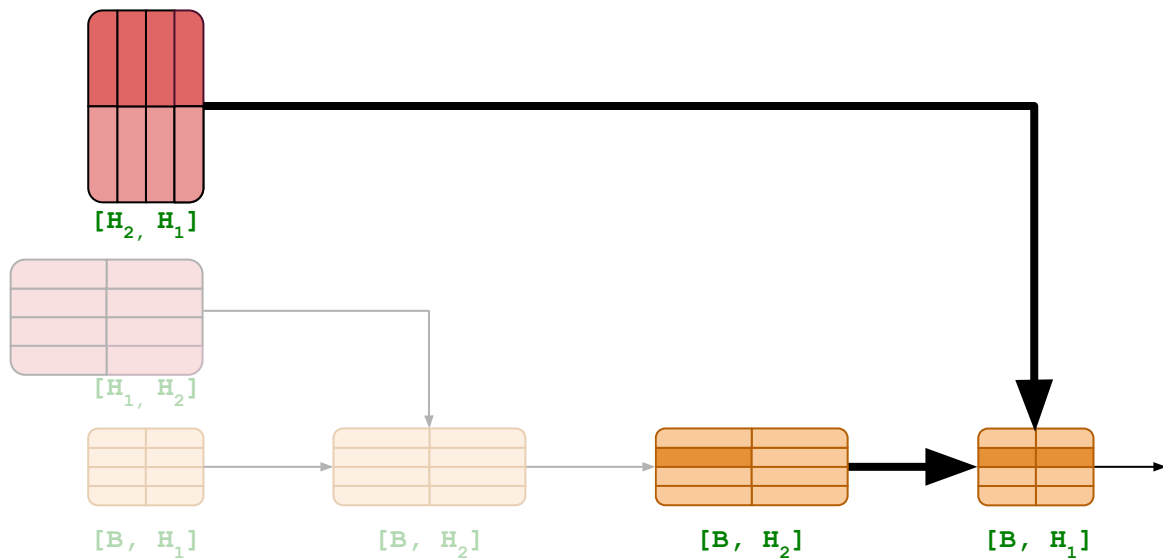


Not really...



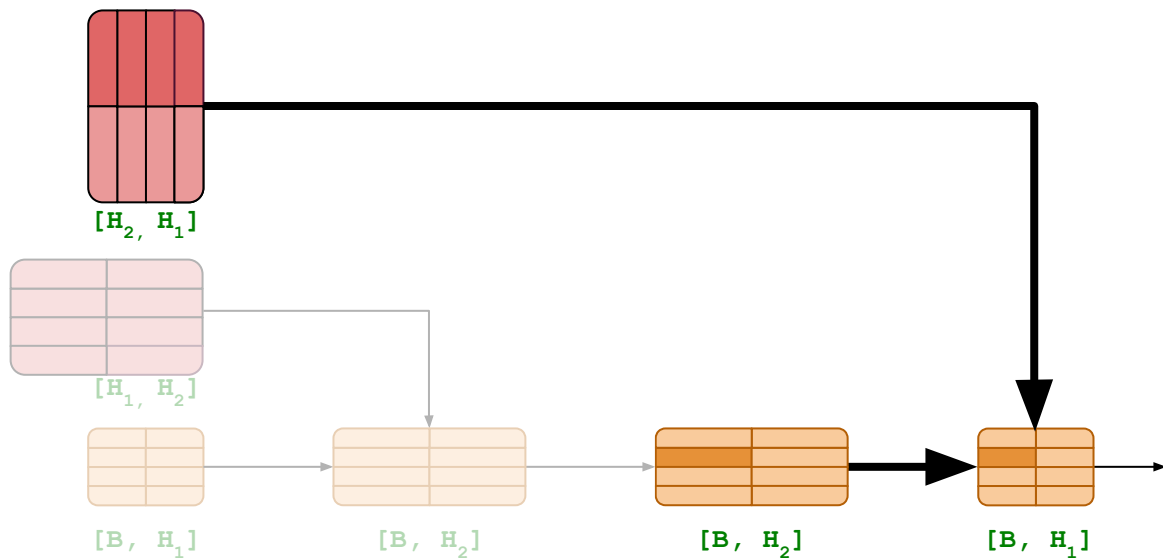
Reduce scatter Z_{yx} on X

$[B/4, H_1]$

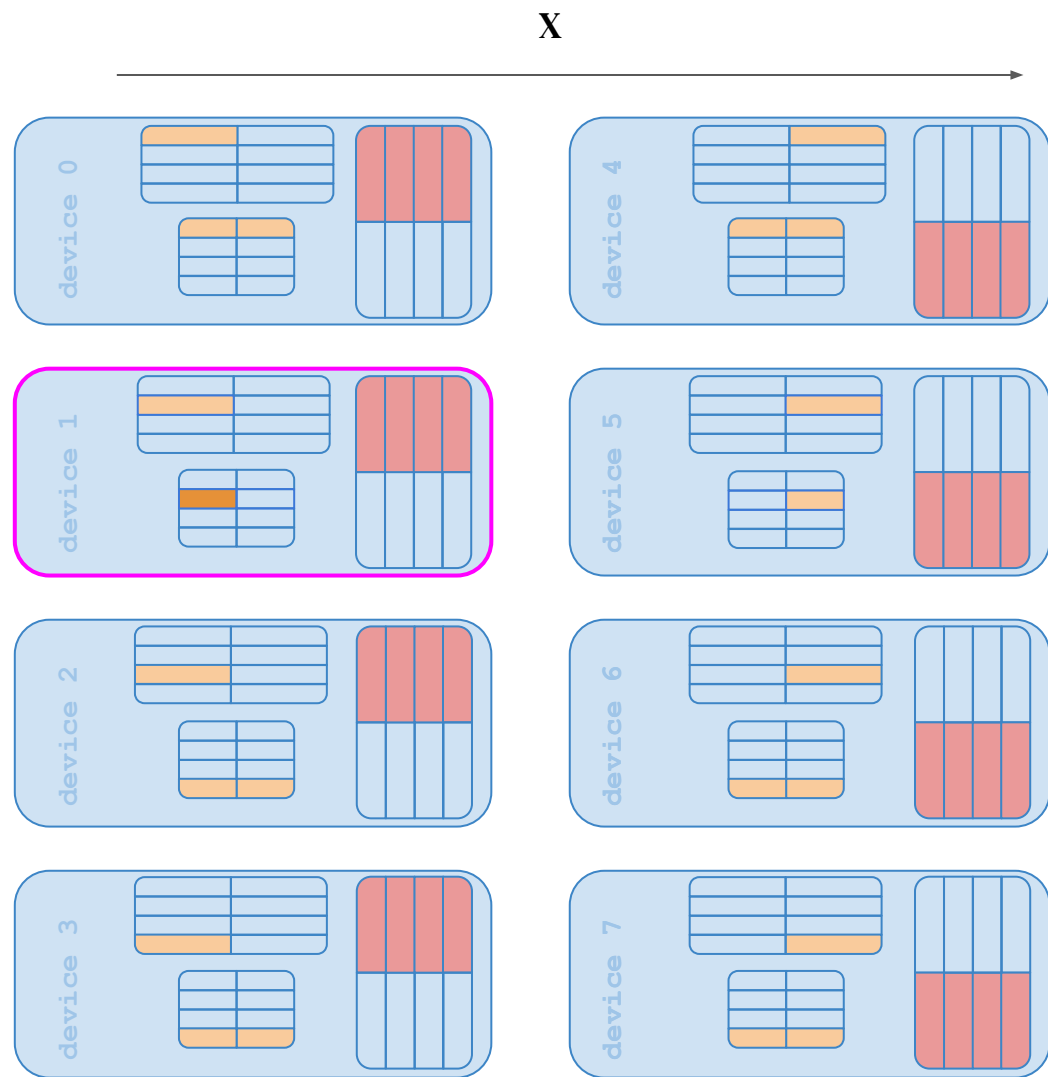


Reduce scatter Z_{yx} on Y

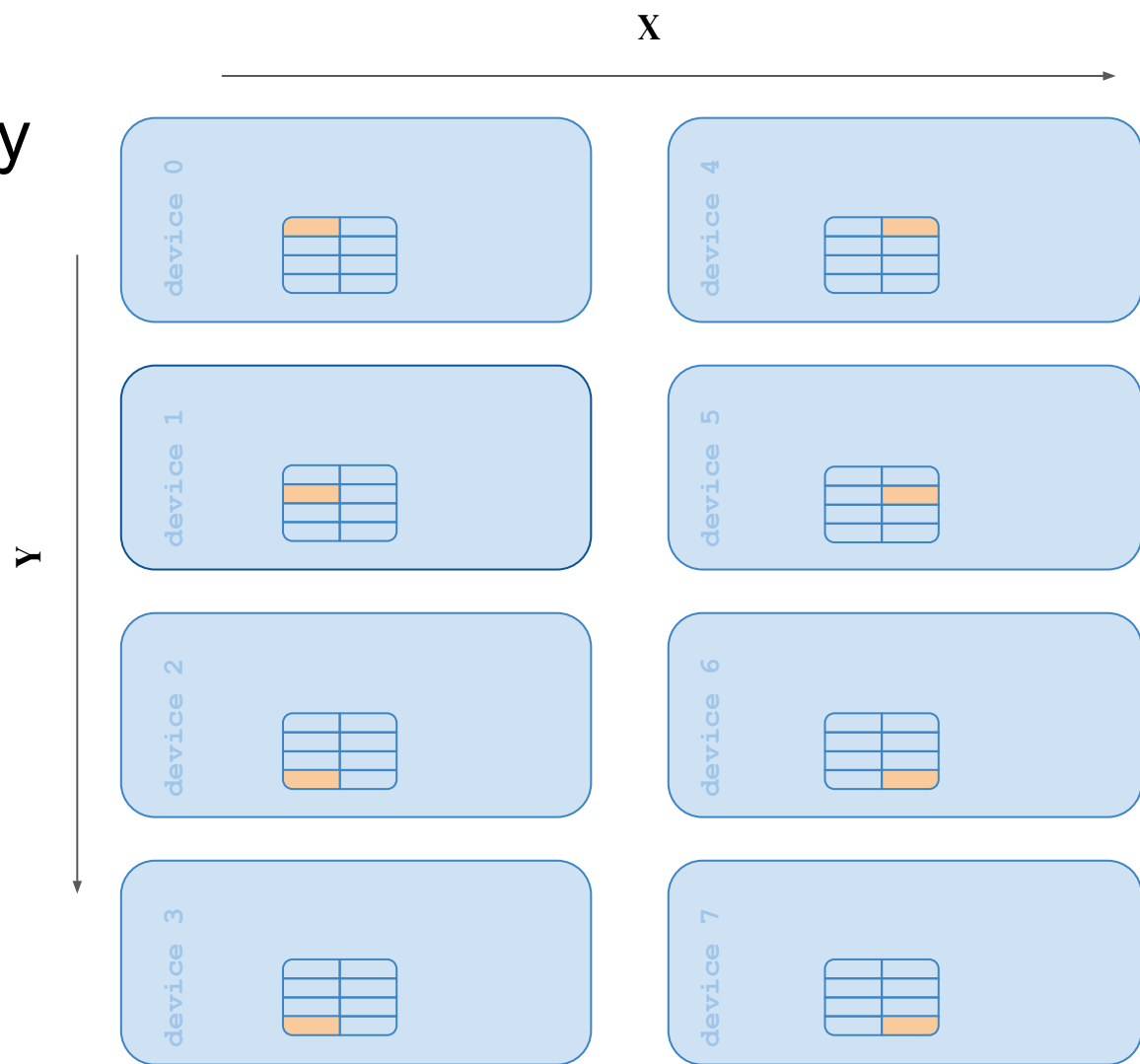
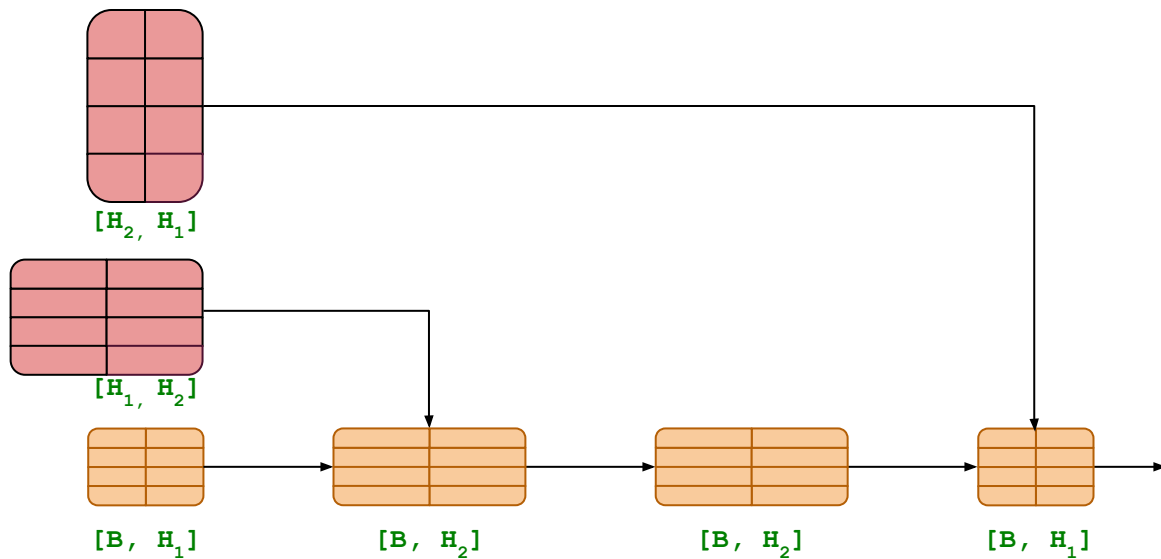
$[B/4, H_1]$



Y



😊 Output is sharded in the same way as input



Summary

- We've described two main approaches to tensor parallelism
 - Megatron sharding
 - 2D parallelism on a mesh
- If you're interested in a deep dive how one would shard a transformer model on a mesh [GSPMD: General and Scalable Parallelization for ML Computation Graphs](#) paper is an excellent introduction
- Specific kind of employed parallelism is dependent on what hardware we're using
 - TPU pods?
 - 2D tensor parallelism works great due to fast interconnect
 - GPU clusters?
 - Combination of pipelining and tensor parallelism

What tools does JAX provide?

- `pmap`
 - Explicit parallel multi-device programming with leading device dimension and manually calling collective ops on reduction axes (JAX documentation [has an excellent introduction](#) to `pmap`) .
- `jax.jit` with `jax.Array`
 - Implicit, compiler-based parallelization system based on user annotation and constraint propagation. Calls GSPMD ([paper](#), [blogpost](#)) under the hood, which rewrites the computation and inserts collective ops.
 - Compiler can will try to be smart, but one can enforce strict constraints with `jax.lax.with_sharding_constraint` ([official JAX tutorial](#))

pmap vs GSPMD

```
A_pmap = shard_on_columns(A, 8)
B_pmap = shard_on_rows(B, 8)

@partial(jax.pmap, axis_name="pmap_axis")
def dot_pmap(x, y):
    return jax.lax.psum(x @ y, 'pmap_axis')
```

```
mesh = jax.sharding.Mesh(
    mesh_utils.create_device_mesh((8, )),
    axis_names=( "x", ))

dot_explicit = jax.jit(
    lambda x, y: jnp.dot(x, y),
    in_shardings=(
        NamedSharding(mesh, PartitionSpec( None, "x")),
        NamedSharding(mesh, PartitionSpec( "x", None))
    )
)
```