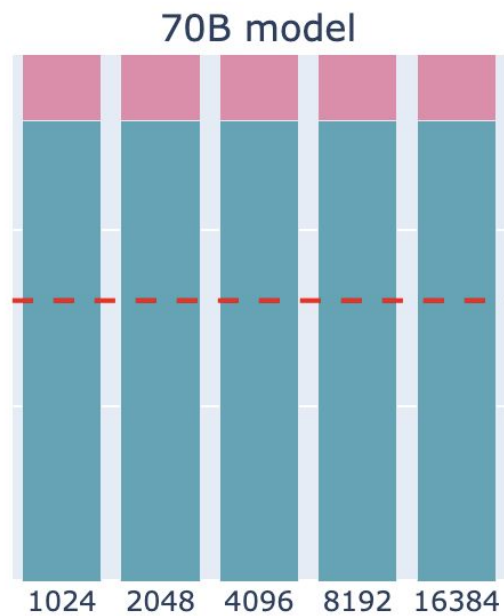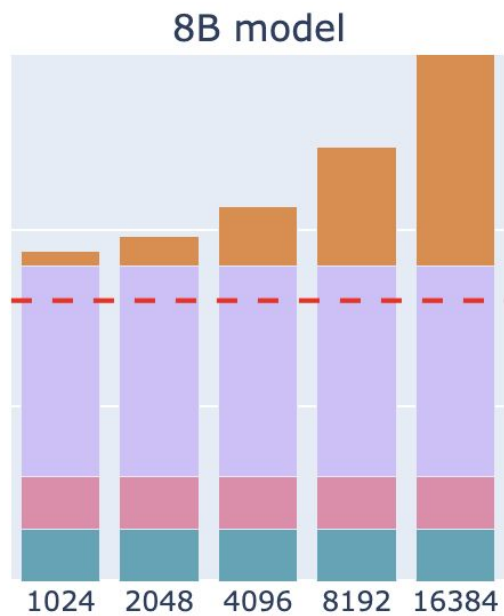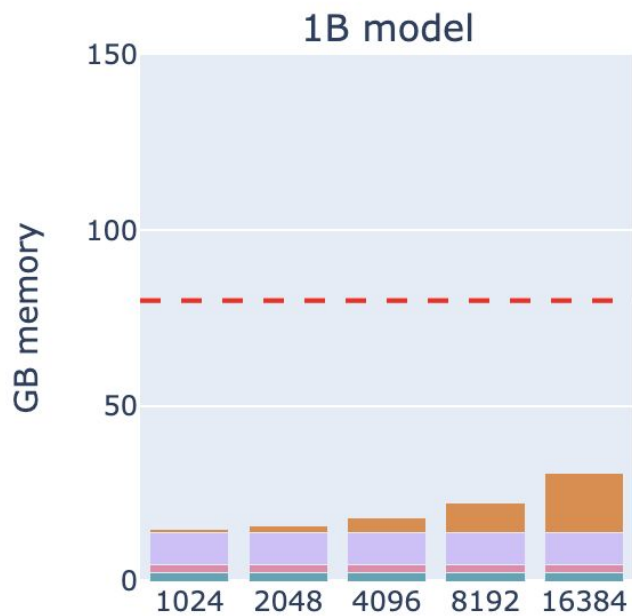# Data Parallelism [ZERO:]

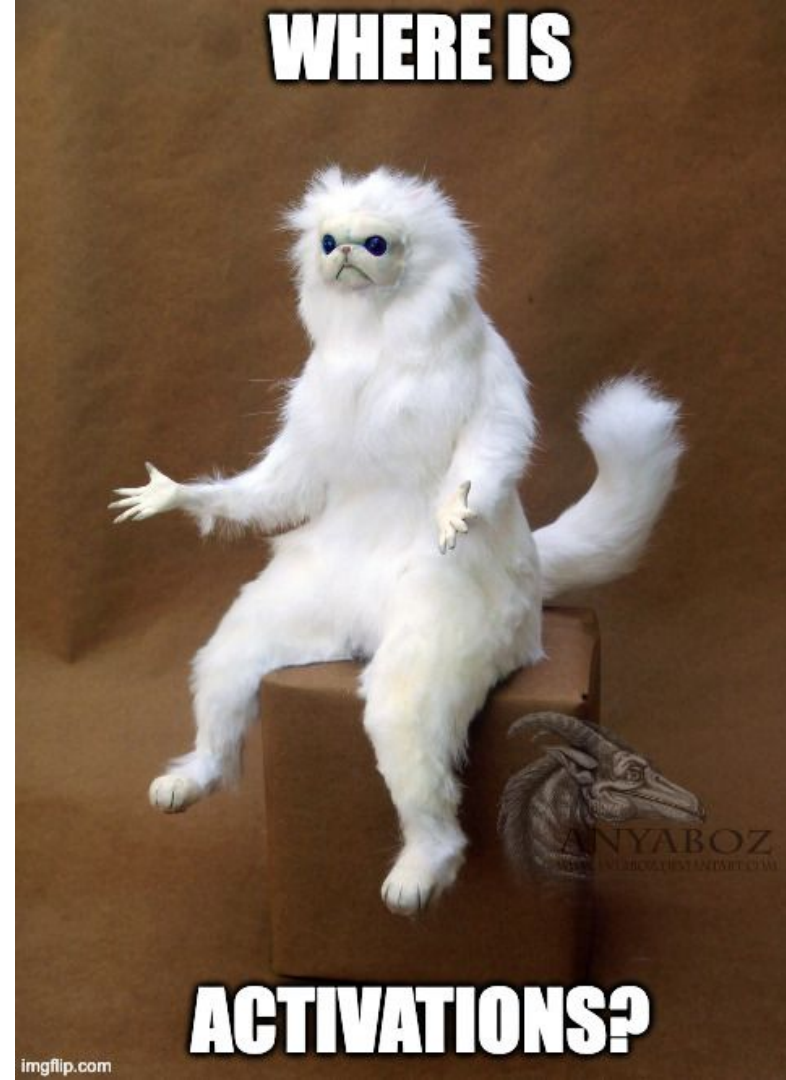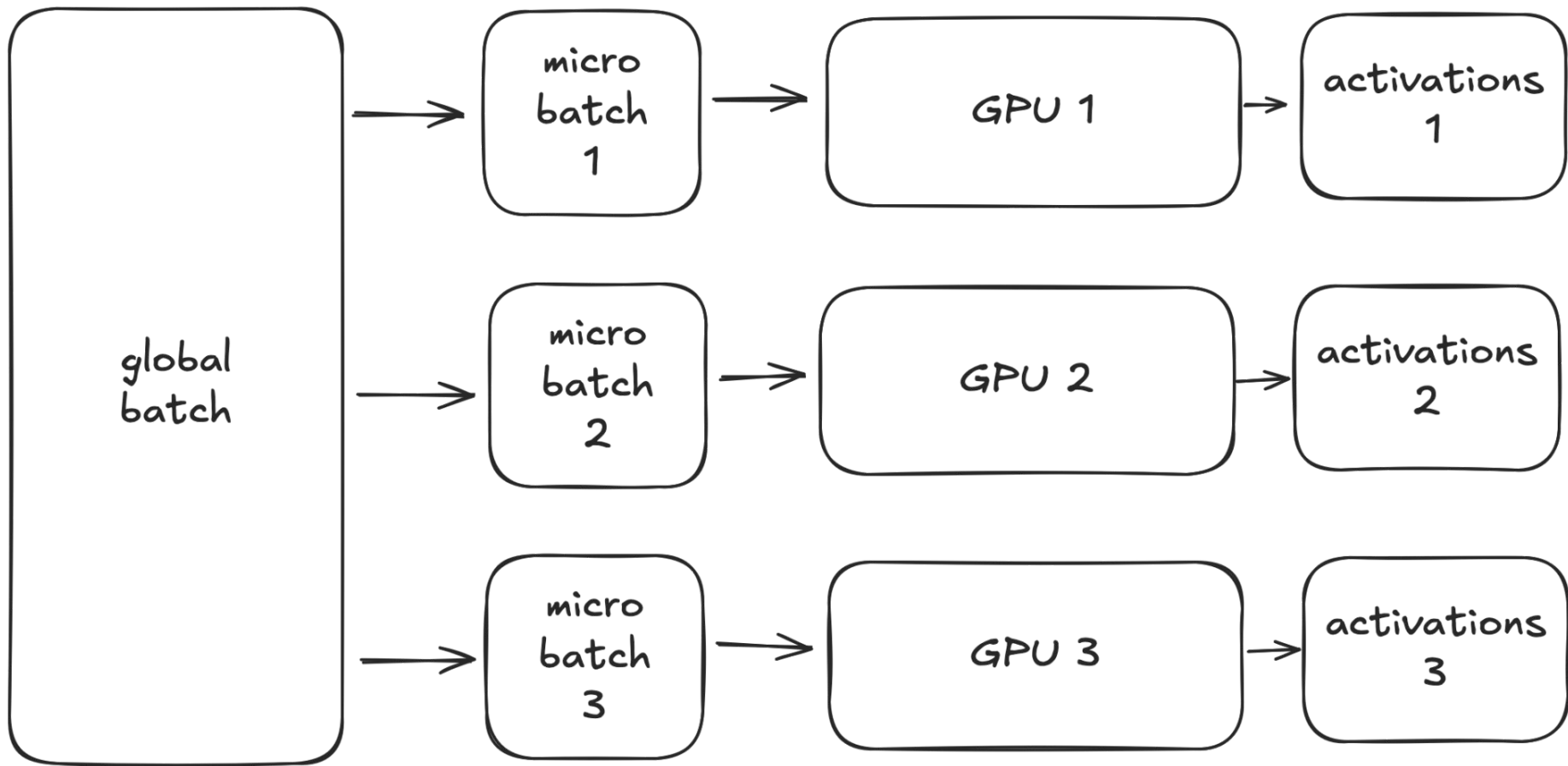made with ❤️ for "Little ML book club"

# Types of ZeRO

1. Optimizer
2. Optimizer + Gradients
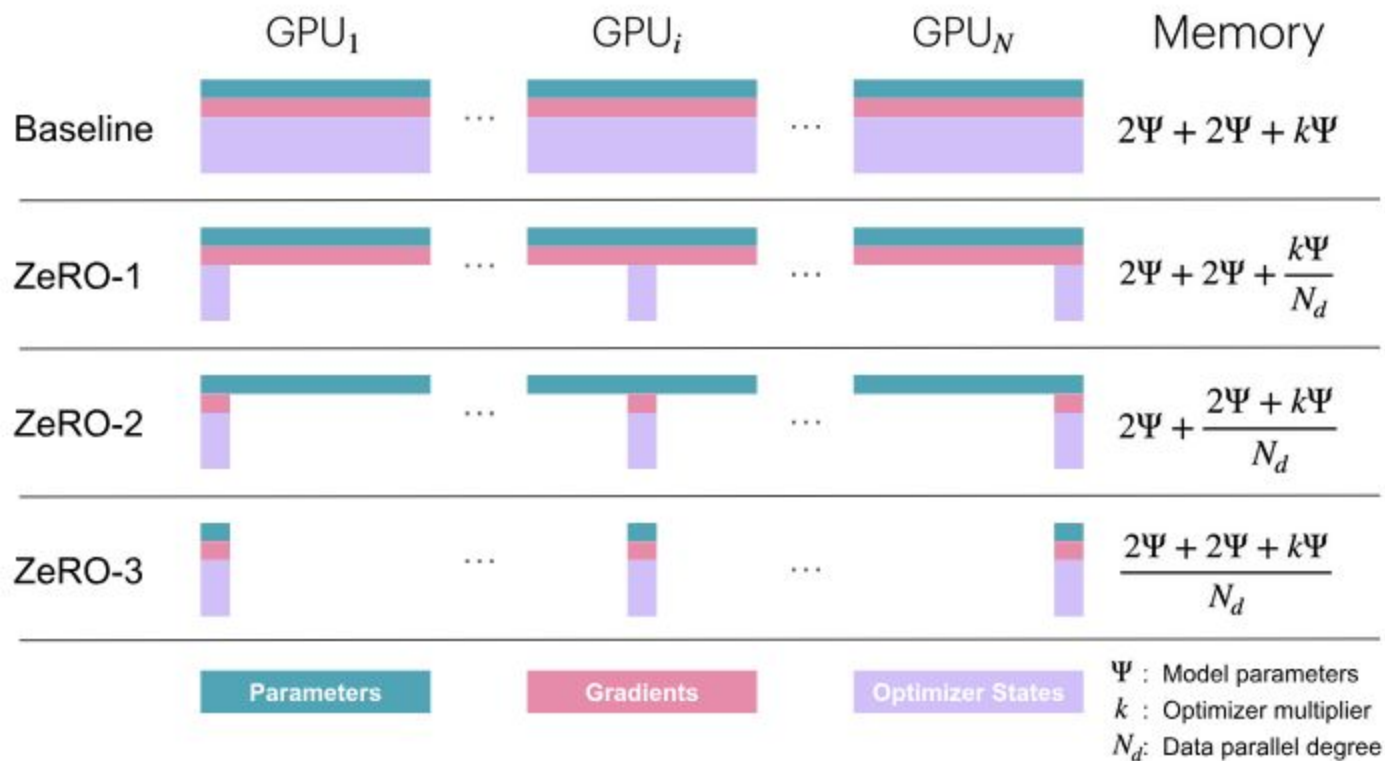3. Optimizer + Gradients + Weights

# Types of ZeRO

1. Optimizer
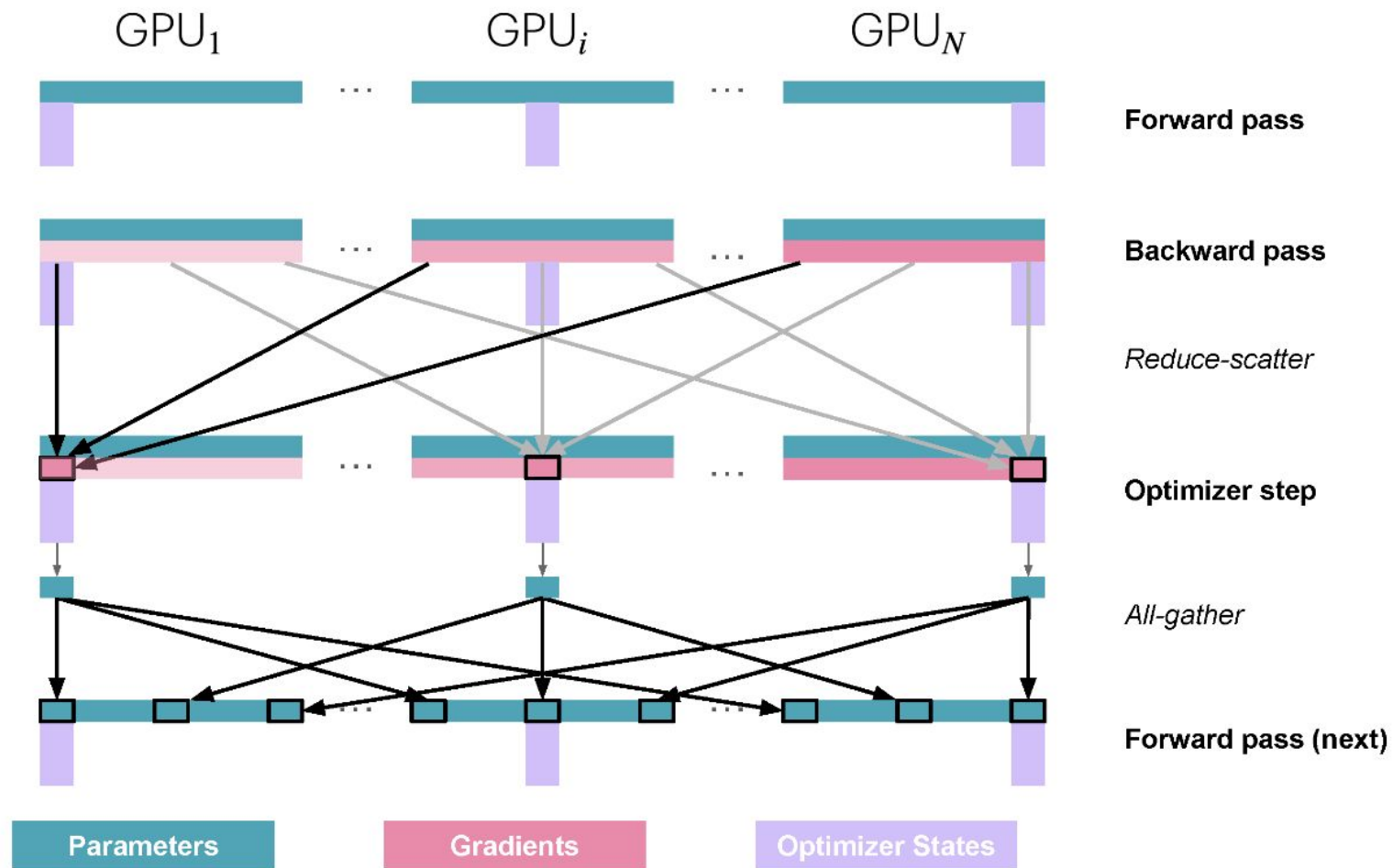2. Optimizer + Gradients
3. Optimizer + Gradients + Weights

- Model's parameters (half precision; i.e., BF16/FP16): $2\Psi$
- Model's gradients (half precision; i.e., BF16/FP16): $2\Psi$
- Model's parameters in FP32 and optimizer states: $4\Psi + (4\Psi + 4\Psi)$
- Model's gradients in FP32: $4\Psi$ (optional, only included if we want to accumulate gradients in FP32)
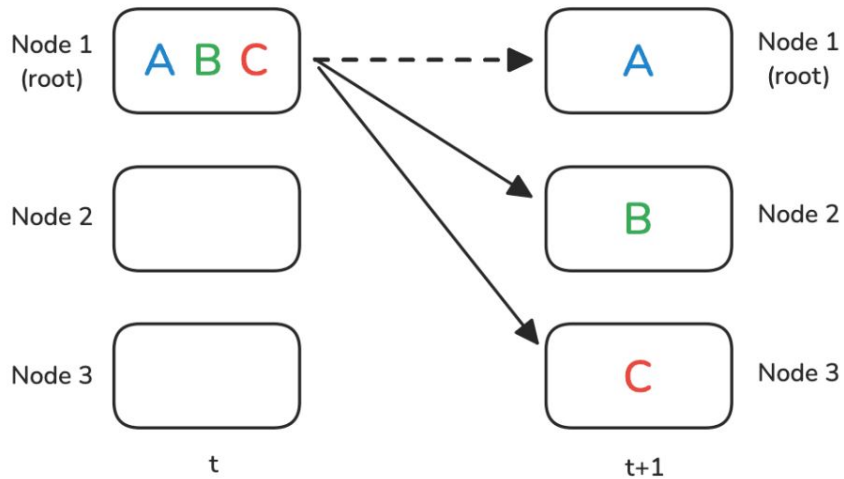
|  | GPU$_1$ | GPU$_i$ | GPU$_N$ | Memory |
|---|---|---|---|---|
| Baseline | | | | $2\Psi + 2\Psi + k\Psi$ |
| ZeRO-1 | | | | $2\Psi + 2\Psi + \dfrac{k\Psi}{N_d}$ |
| ZeRO-2 | | | | $2\Psi + \dfrac{2\Psi + k\Psi}{N_d}$ |
| ZeRO-3 | | | | $\dfrac{2\Psi + 2\Psi + k\Psi}{N_d}$ |

Parameters    Gradients    Optimizer States

$\Psi$ : Model parameters
$k$ : Optimizer multiplier
$N_d$: Data parallel degree

# ZeRO 1

| GPU$_1$ | GPU$_i$ | GPU$_N$ | |
|---|---|---|---|

Forward pass

Backward pass

*Reduce-scatter*

Optimizer step

*All-gather*

Forward pass (next)

Parameters    Gradients    Optimizer States

# Scatter

Node 1 (root): A B C

Node 1 (root): A
Node 2: B
Node 3: C

t → t+1

# ReduceScatter

Node 1 (root): A
Node 2: B
Node 3: C

f()

Node 1 (root): x
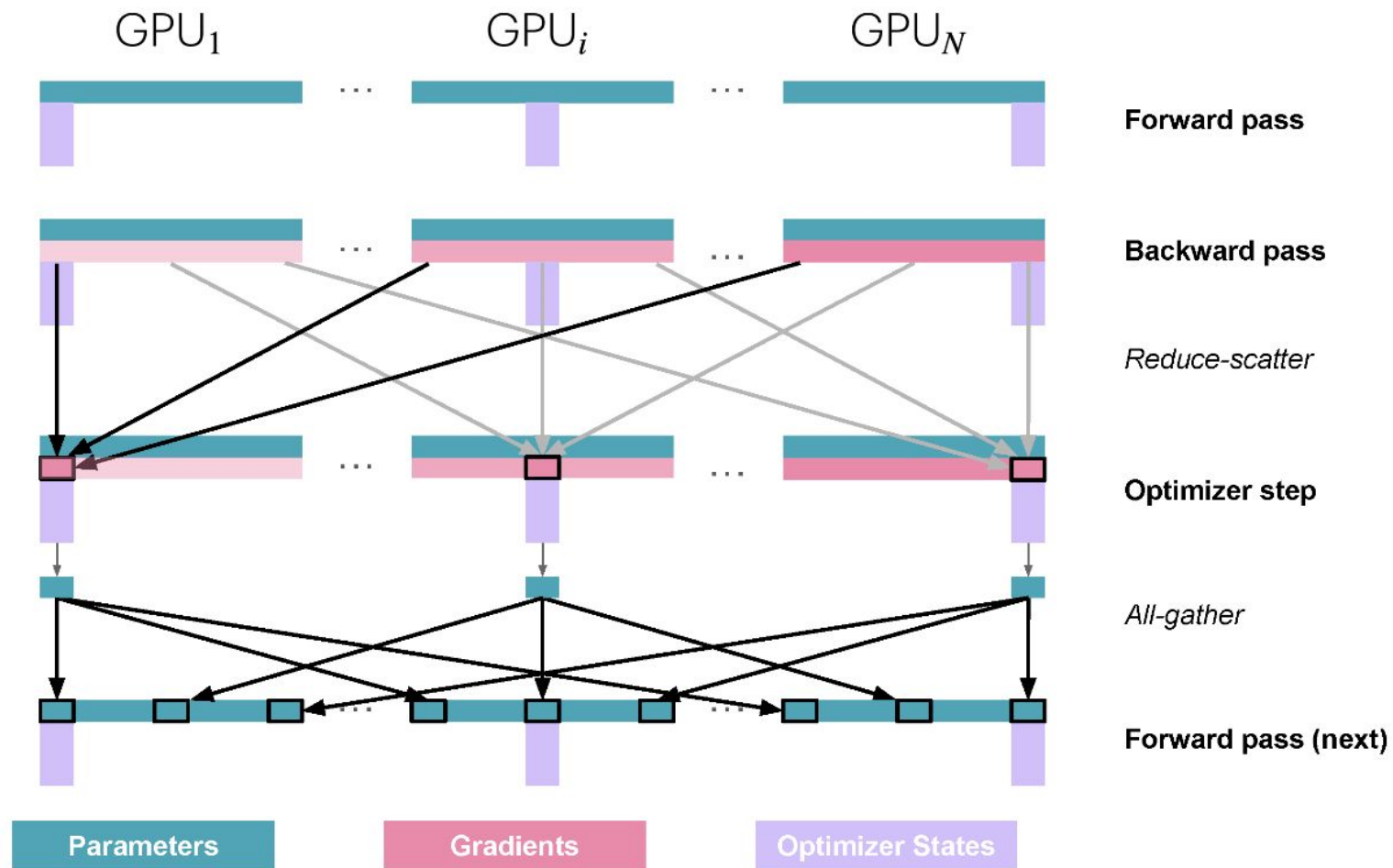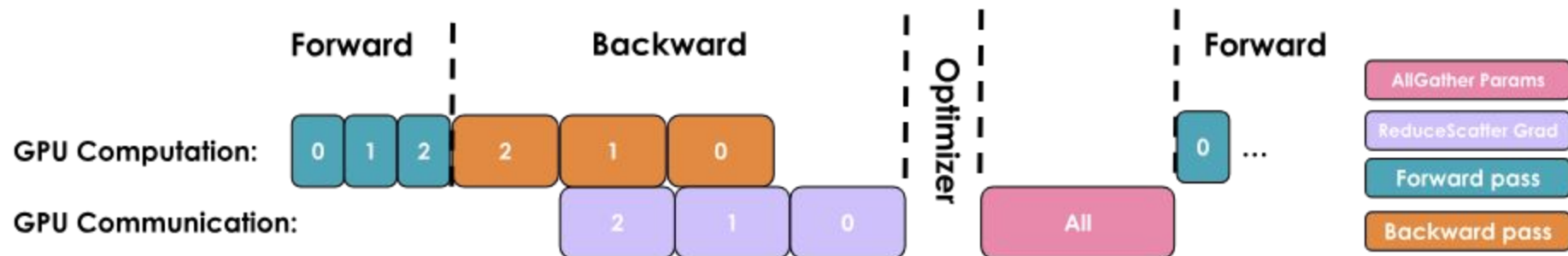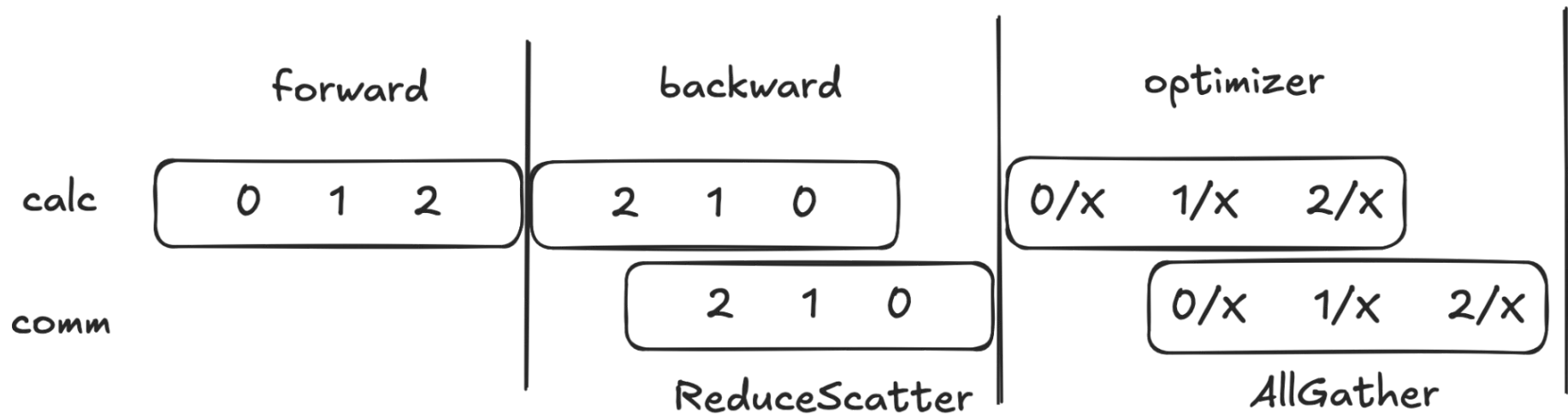Node 2: y
Node 3: z

t → t+1

```python
def example_reduce_scatter():
    rank = dist.get_rank()
    world_size = dist.get_world_size()
    input_tensor = [
        torch.tensor([(rank + 1) * i for i in range(1, 3)], dtype=torch.float32).cuda()**(j+1)
        for j in range(world_size)
        ]
    output_tensor = torch.zeros(2, dtype=torch.float32).cuda()
    dist.reduce_scatter(output_tensor, input_tensor, op=dist.ReduceOp.SUM)
```

```
Before ReduceScatter on rank 0: [tensor([1., 2.], device='cuda:0'),
                                 tensor([1., 4.], device='cuda:0'),
                                 tensor([1., 8.], device='cuda:0')]
Before ReduceScatter on rank 1: [tensor([2., 4.], device='cuda:1'),
                                 tensor([4., 16.], device='cuda:1'),
                                 tensor([8., 64.], device='cuda:1')]
Before ReduceScatter on rank 2: [tensor([3., 6.], device='cuda:2'),
                                 tensor([9., 36.], device='cuda:2'),
                                 tensor([27., 216.], device='cuda:2')]

After ReduceScatter on rank 0: tensor([6., 12.], device='cuda:0')
After ReduceScatter on rank 1: tensor([14., 56.], device='cuda:1')
After ReduceScatter on rank 2: tensor([36., 288.], device='cuda:2')
```
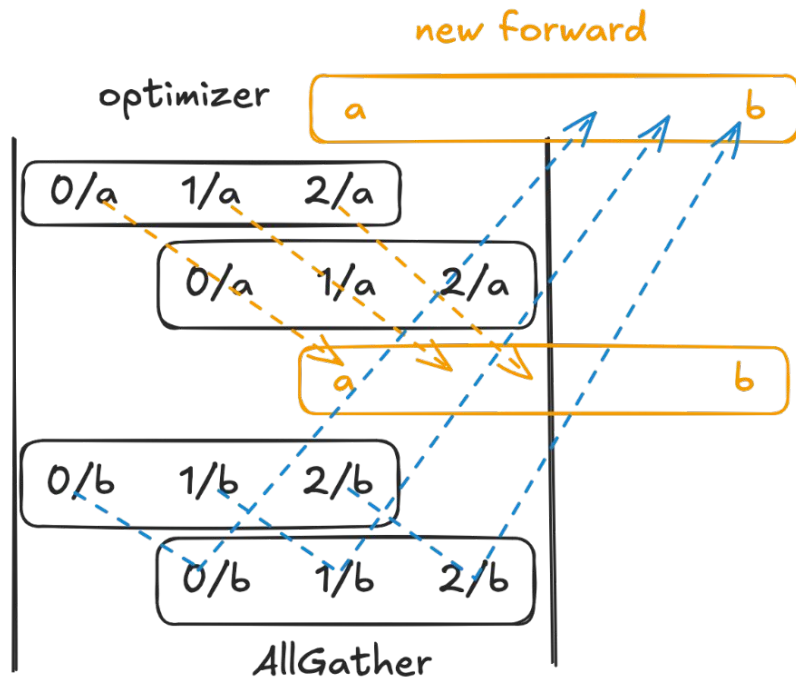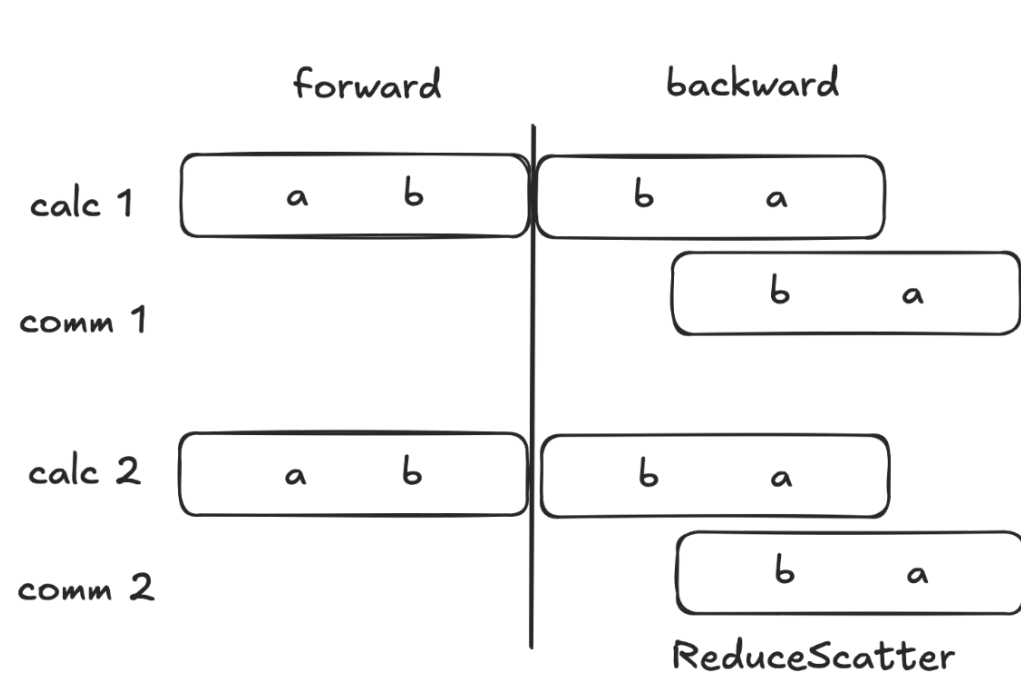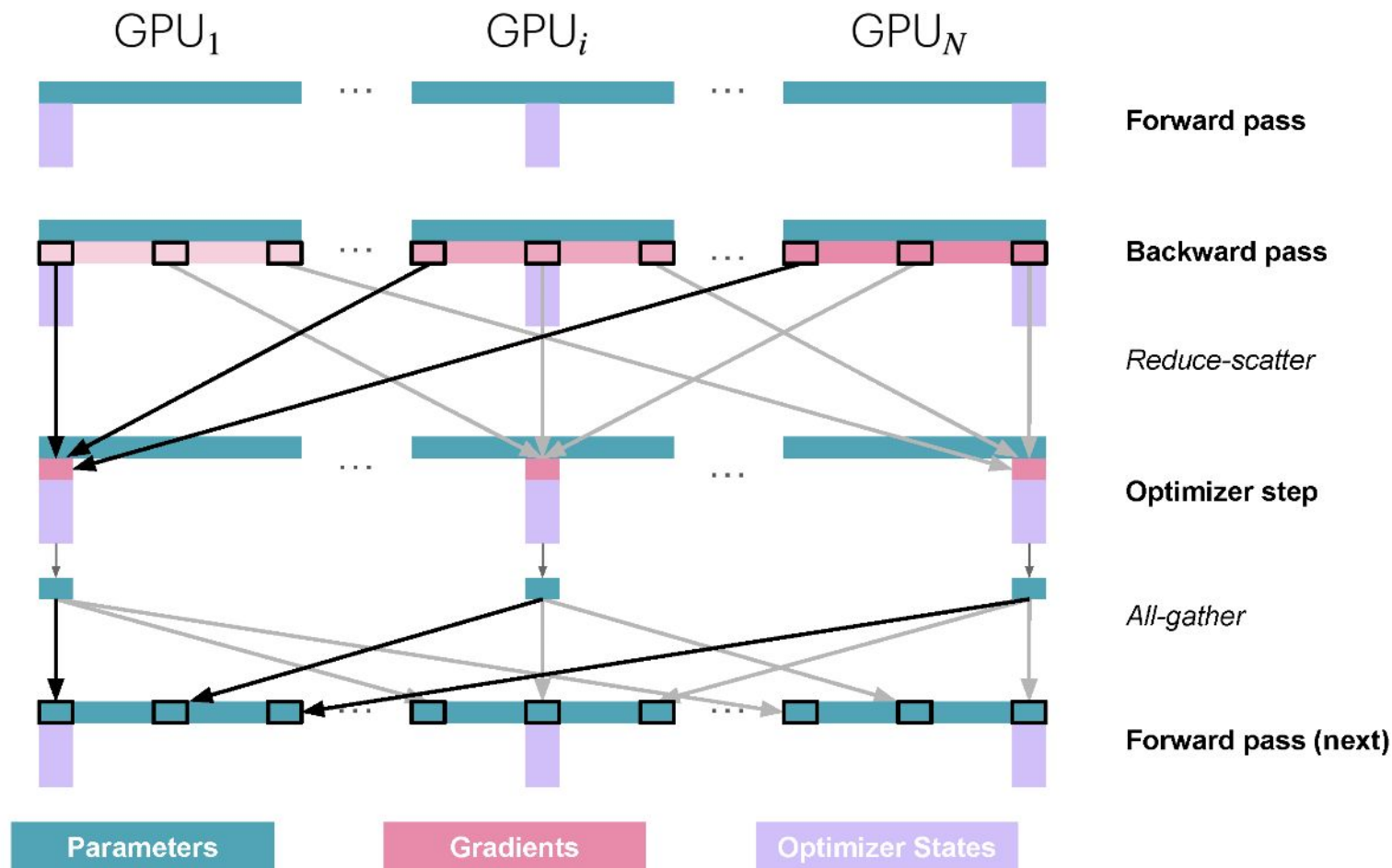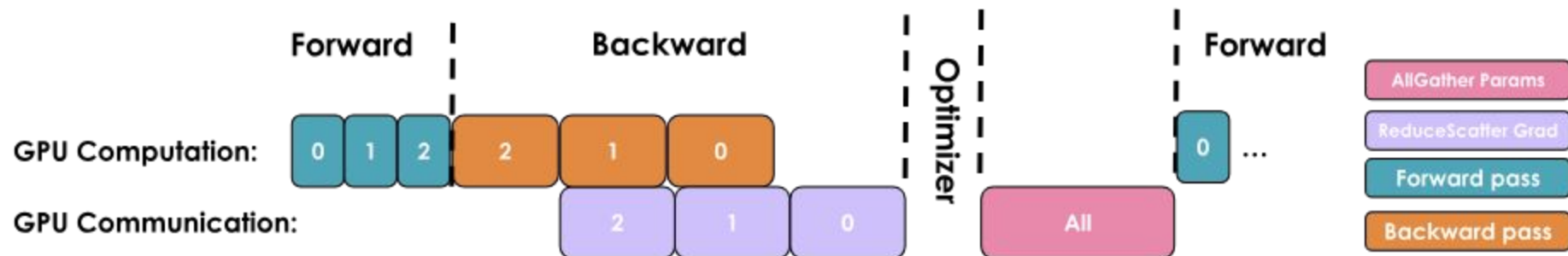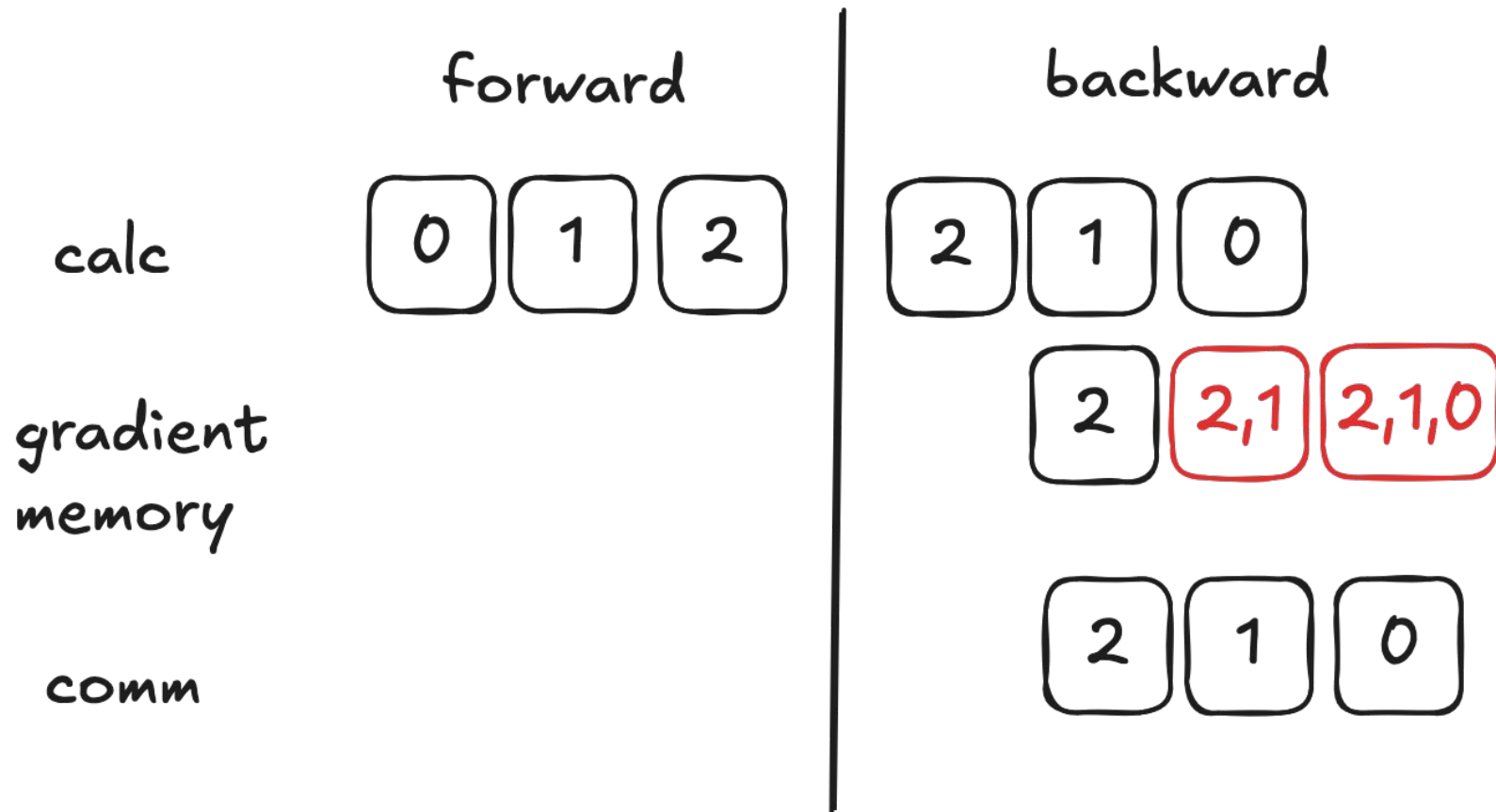
GPU$_1$     ⋯     GPU$_i$     ⋯     GPU$_N$

**Forward pass**

**Backward pass**

*Reduce-scatter*

**Optimizer step**

*All-gather*

**Forward pass (next)**

**Parameters**     **Gradients**     **Optimizer States**

# ZeRO 2

GPU$_1$ ... GPU$_i$ ... GPU$_N$

Forward pass

Backward pass

Reduce-scatter

Optimizer step

All-gather

Forward pass (next)

Parameters  Gradients  Optimizer States

# ZeRO 3

FSDP (Fully Sharded Data Parallelism) in pytroch

GPU$_1$ ⋯ GPU$_i$ ⋯ GPU$_N$

Layer n-1

Reduce-scatter

Layer n

Backward pass

All-gather

Layer n+1

Parameters   Gradients

Forward | Backward

GPU Computation:

GPU Communication:

AllGather Params
ReduceScatter Grads
Forward pass
Backward pass
Parameter Free

forward

gpu 0 calc
gpu 1 calc
gpu 2 calc
gpu 0 IO (out)
gpu 1 IO (out)
gpu 2 IO (out)
gpu 0 IO (in)
gpu 1 IO (in)
gpu 2 IO (in)

Memory Usage for 8B Model

mom,
but I want to play
with veeeeery long
sequences

see you next week for
tensor parallelism