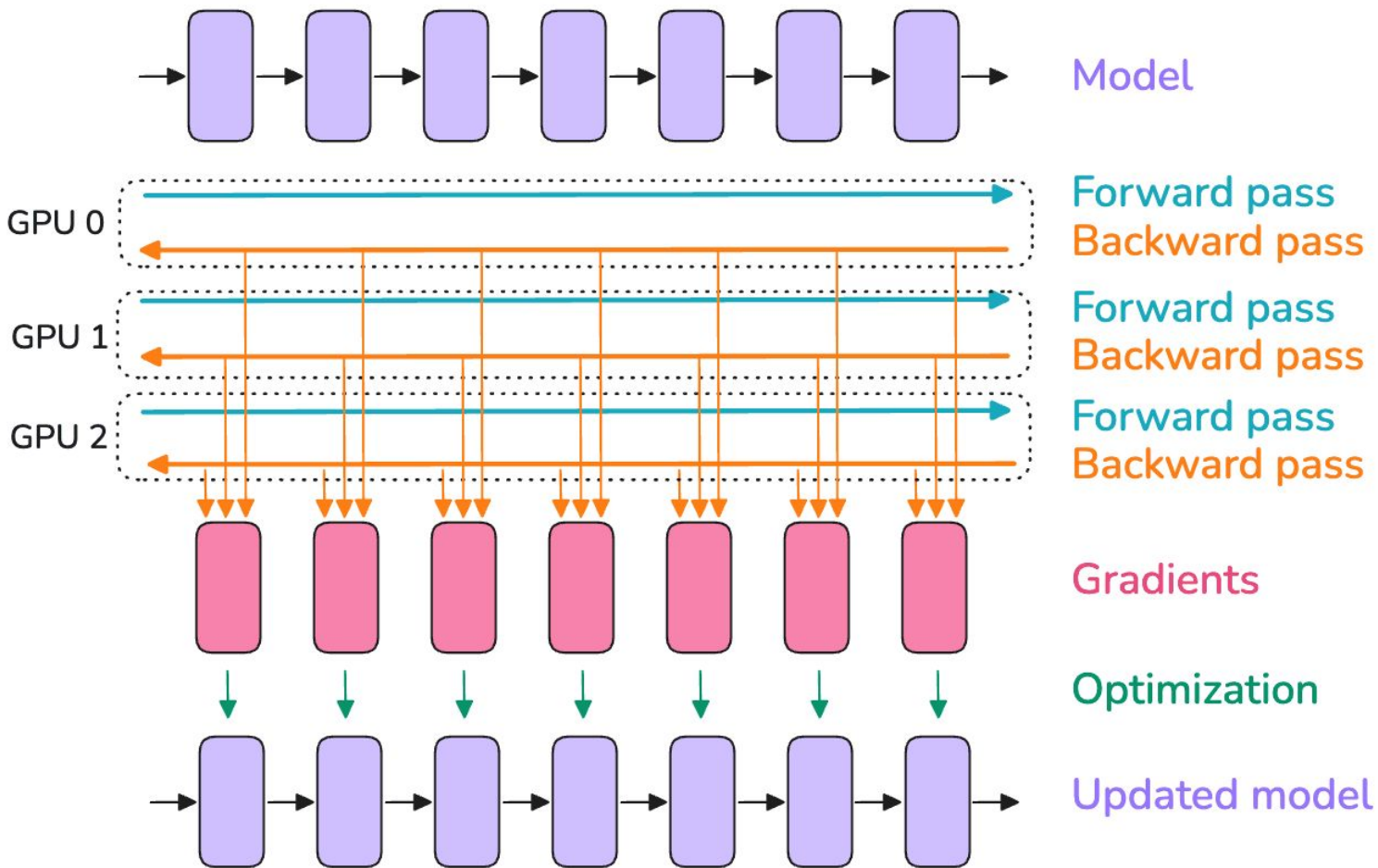
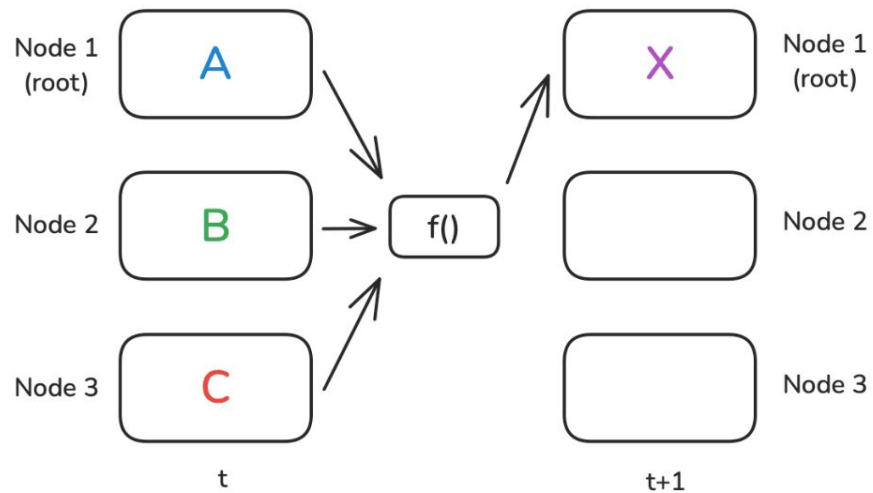


Data Parallelism [:ZERO]

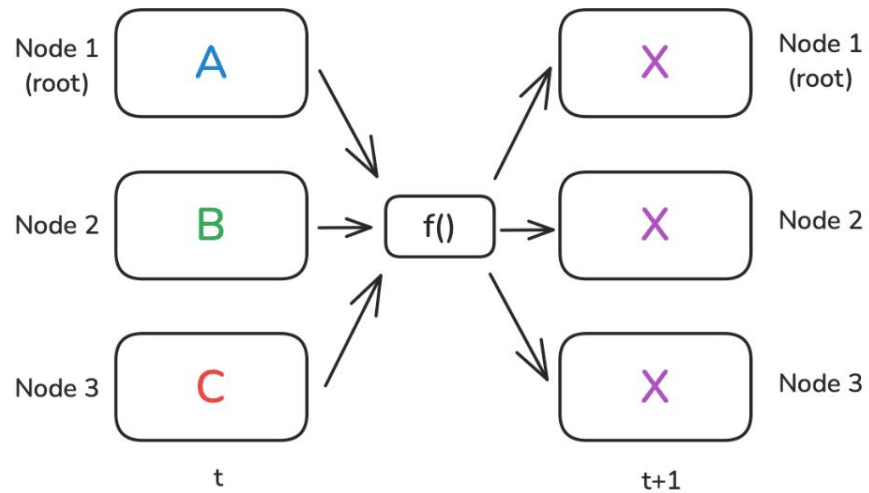
made with  for “Little ML book club”

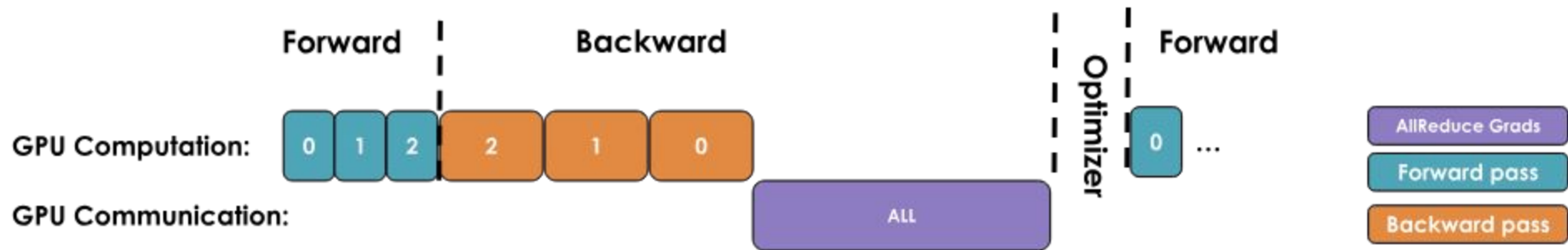


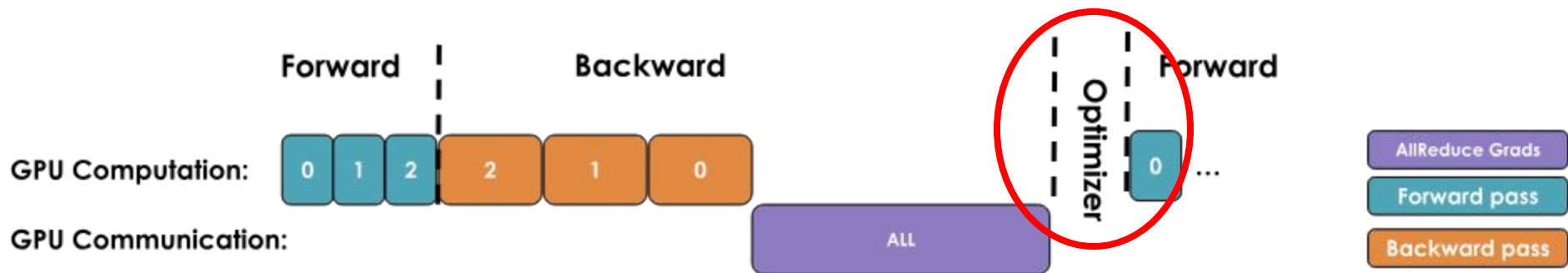
Reduce

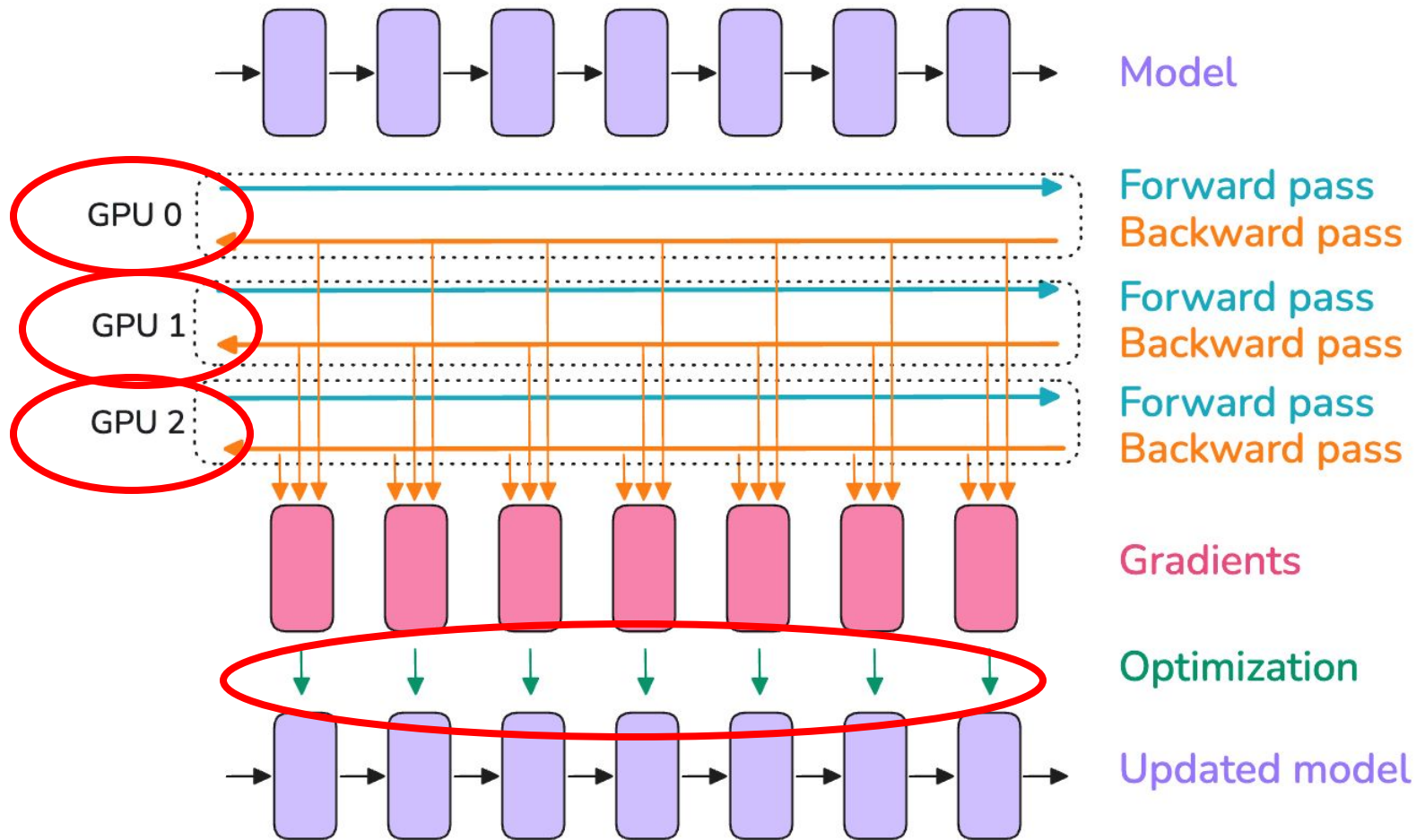


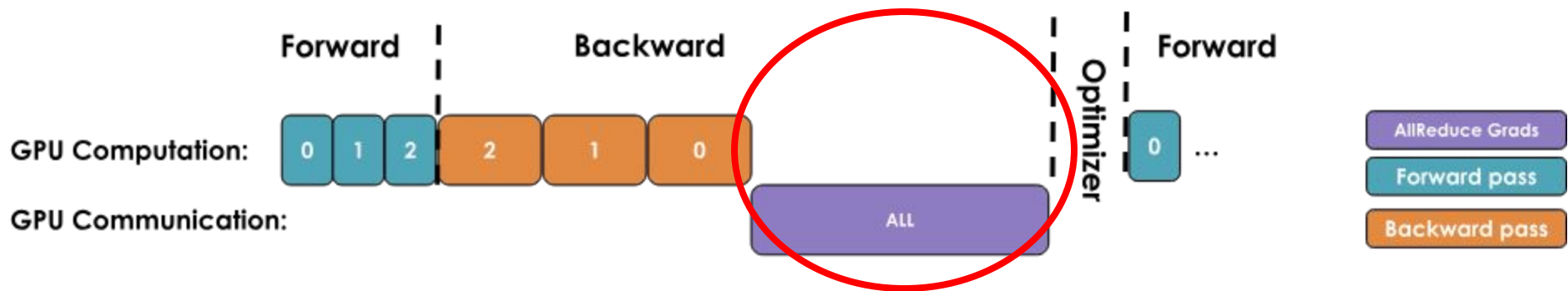
AllReduce

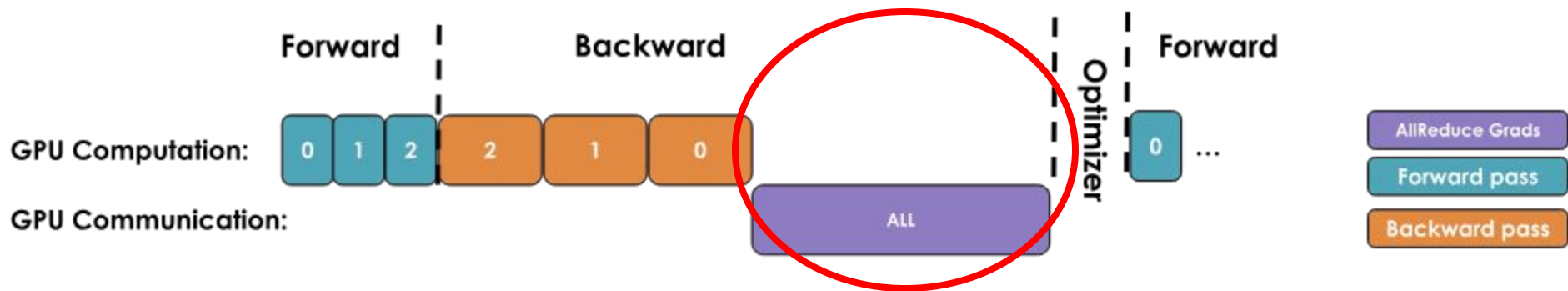




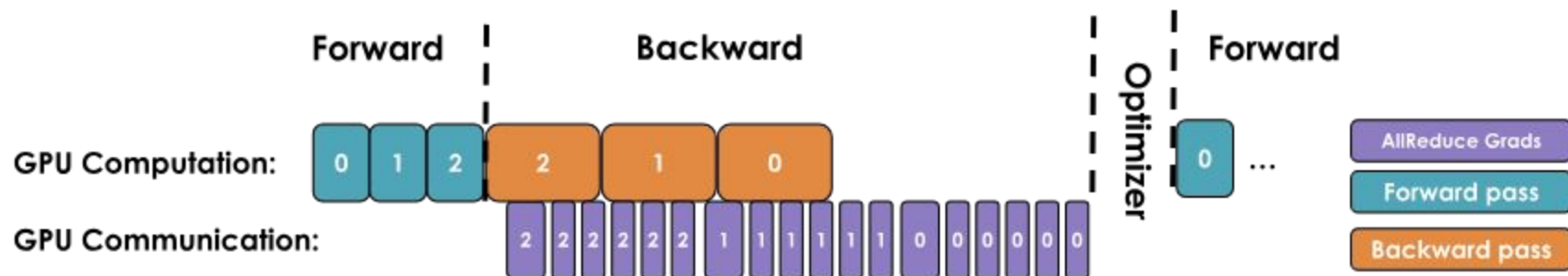


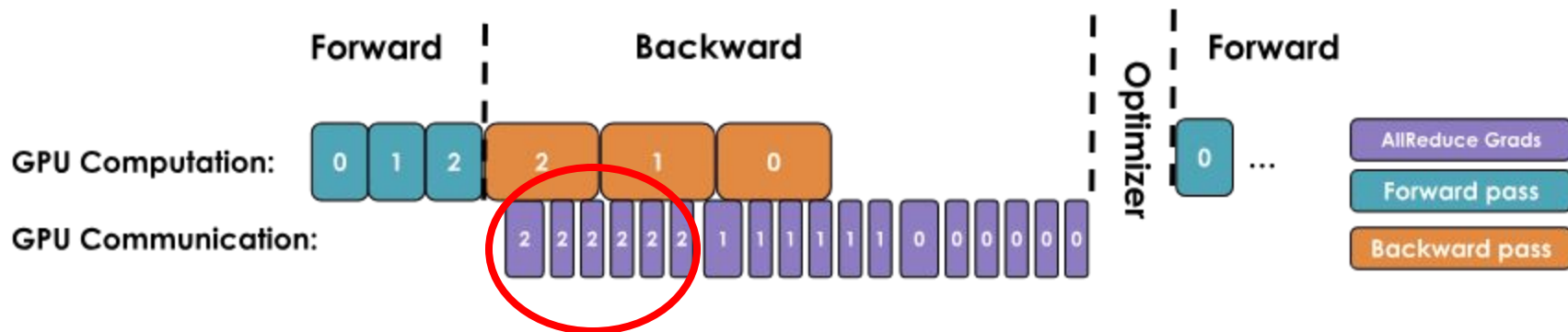


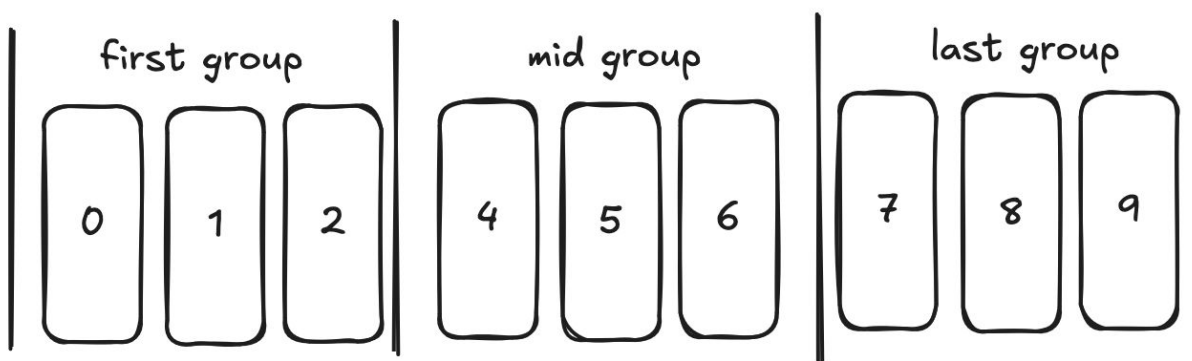
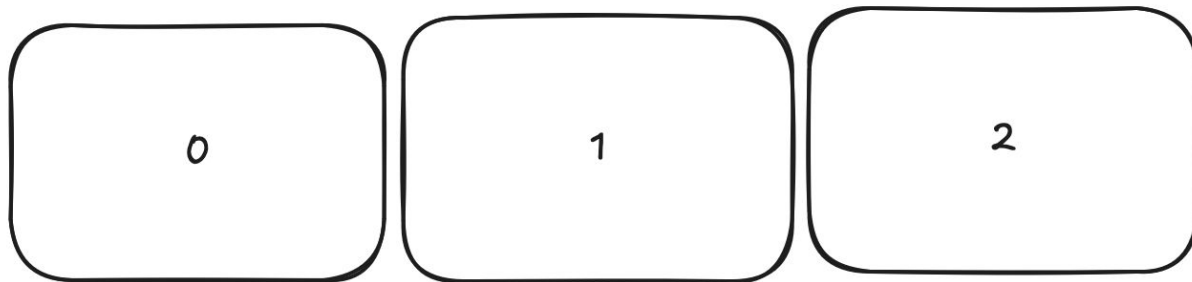




NOTE: Each GPU needs memory for all grads!

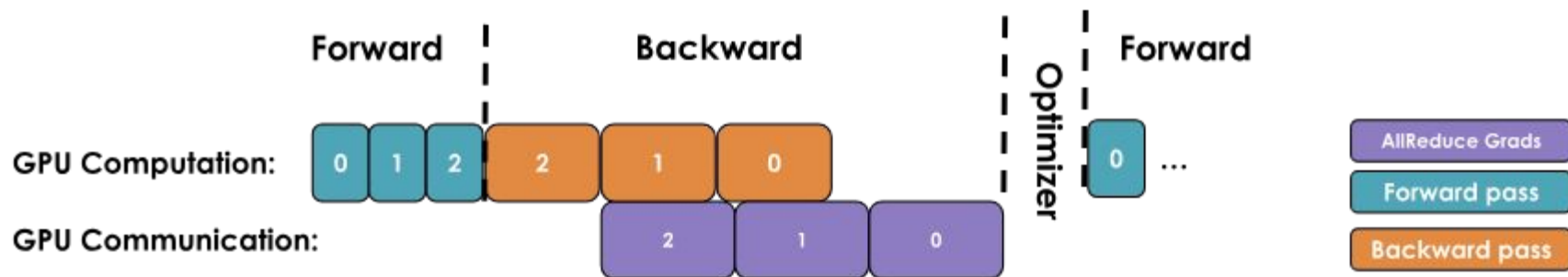
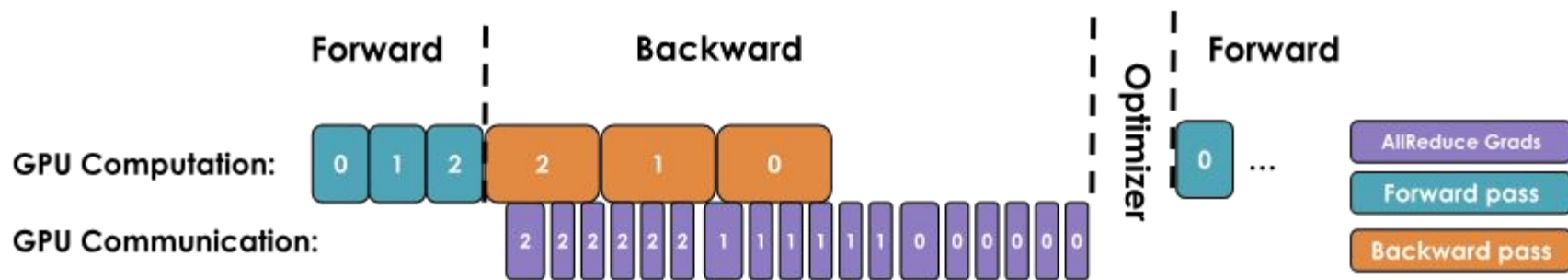


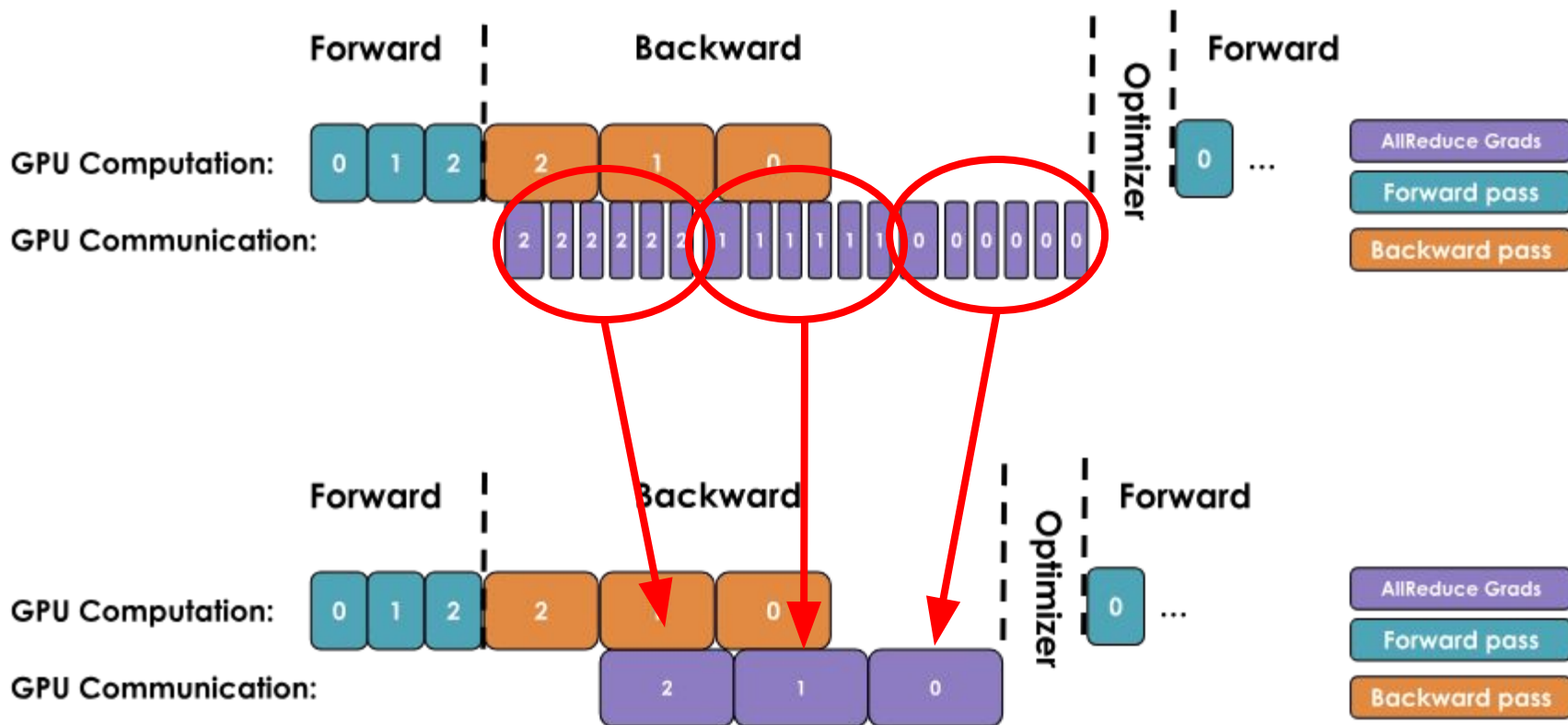




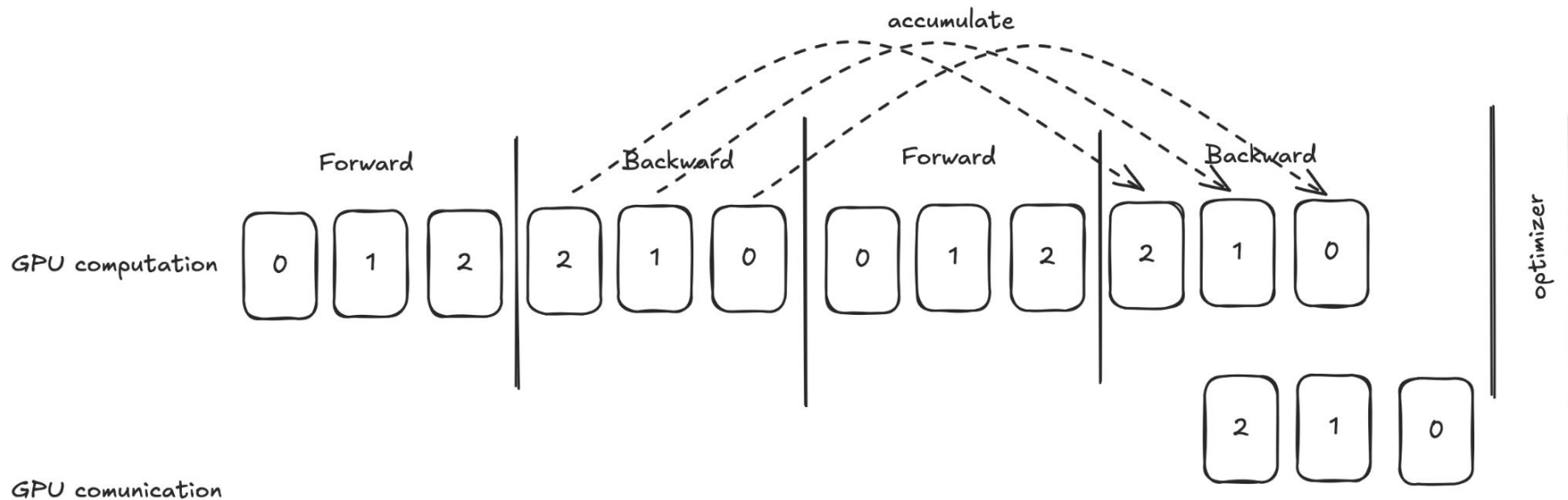


```
def _allreduce_grads(self, grad):  
    """  
    Performs an all-reduce operation to synchronize gradients across multiple processes.  
    """  
    # No synchronization needed during gradient accumulation, except at the final accumulation step.  
    if self.require_backward_grad_sync:  
        dist.all_reduce(  
            grad, op=dist.ReduceOp.SUM, group=pgm.process_group_manager.cp_dp_group  
        )  
        grad /= pgm.process_group_manager.cp_dp_world_size  
    return grad
```





mom,
what if my batch size is
not big enough
for optimal
convergence?

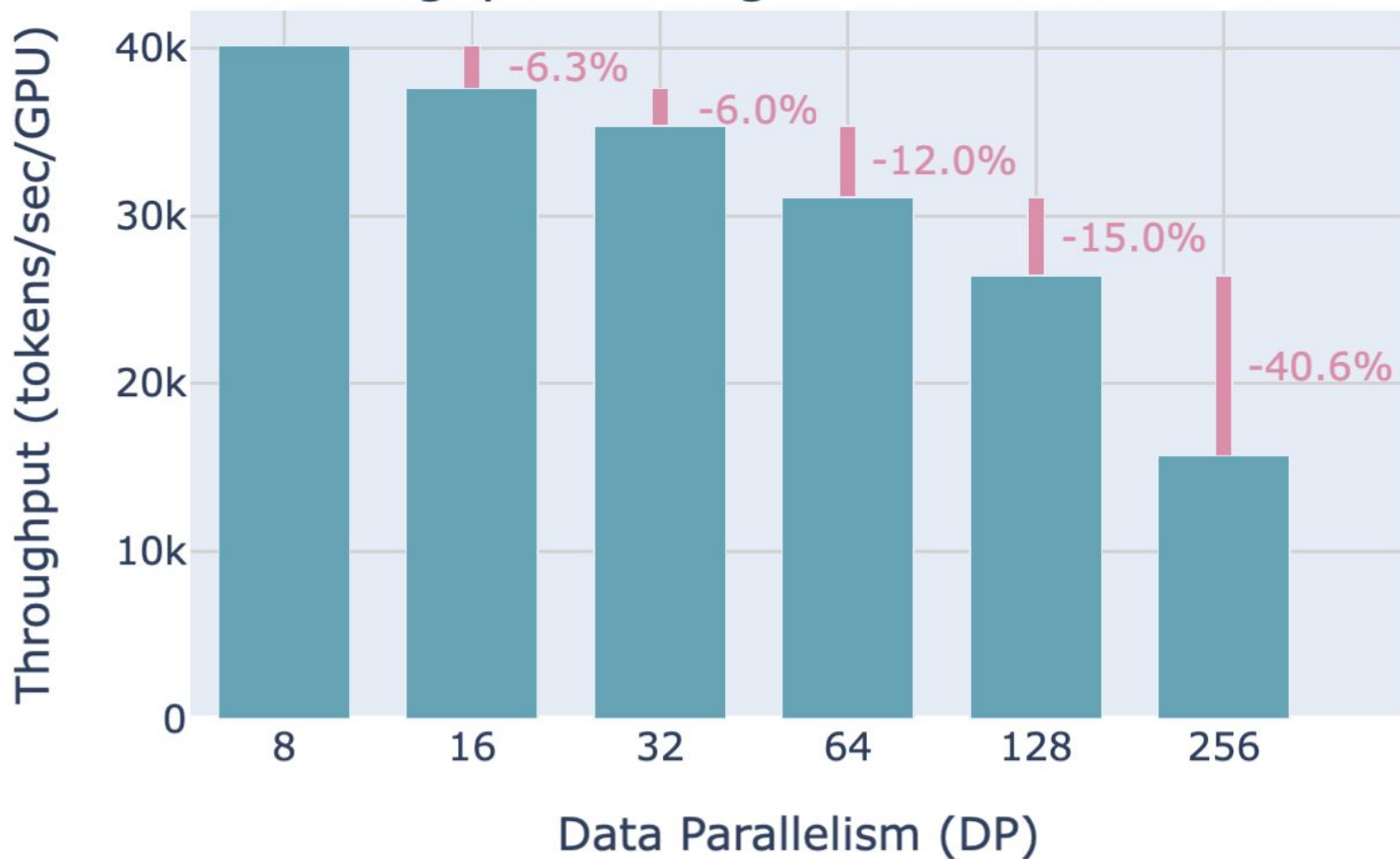


$$bs = gbs = mbs \times grad_acc \times dp$$

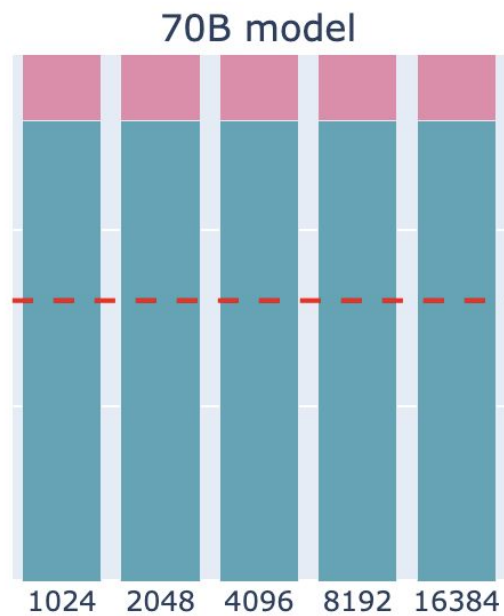
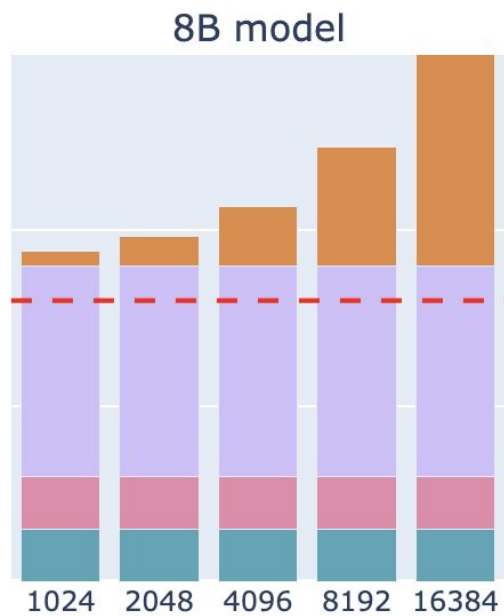
Recipe

1. Find best (global) batch size in tokens (sacred empirical knowledge)
2. Find a sequence length (sacred empirical knowledge). Generally, 2-8k tokens.
3. Find the maximum local batch size on a single GPU by increasing the local batch size until we run out of memory.
4. Find the number of available GPUs for our target (dp). The ratio of gbs to dp gives us the remaining number of gradient accumulation steps needed for the desired gbs .

Throughput Scaling with Data Parallelism



mom,
what if my model is too
large?



see you next week
to solve it