

# 分布式训练系列

## 集合通信原语



BUILDING A BETTER CONNECTED WORLD

Ascend & MindSpore

[www.hiascend.com](http://www.hiascend.com)  
[www.mindspore.cn](http://www.mindspore.cn)

# 关于本内容

## 1. 内容背景

- AI集群+大模型+分布式训练系统

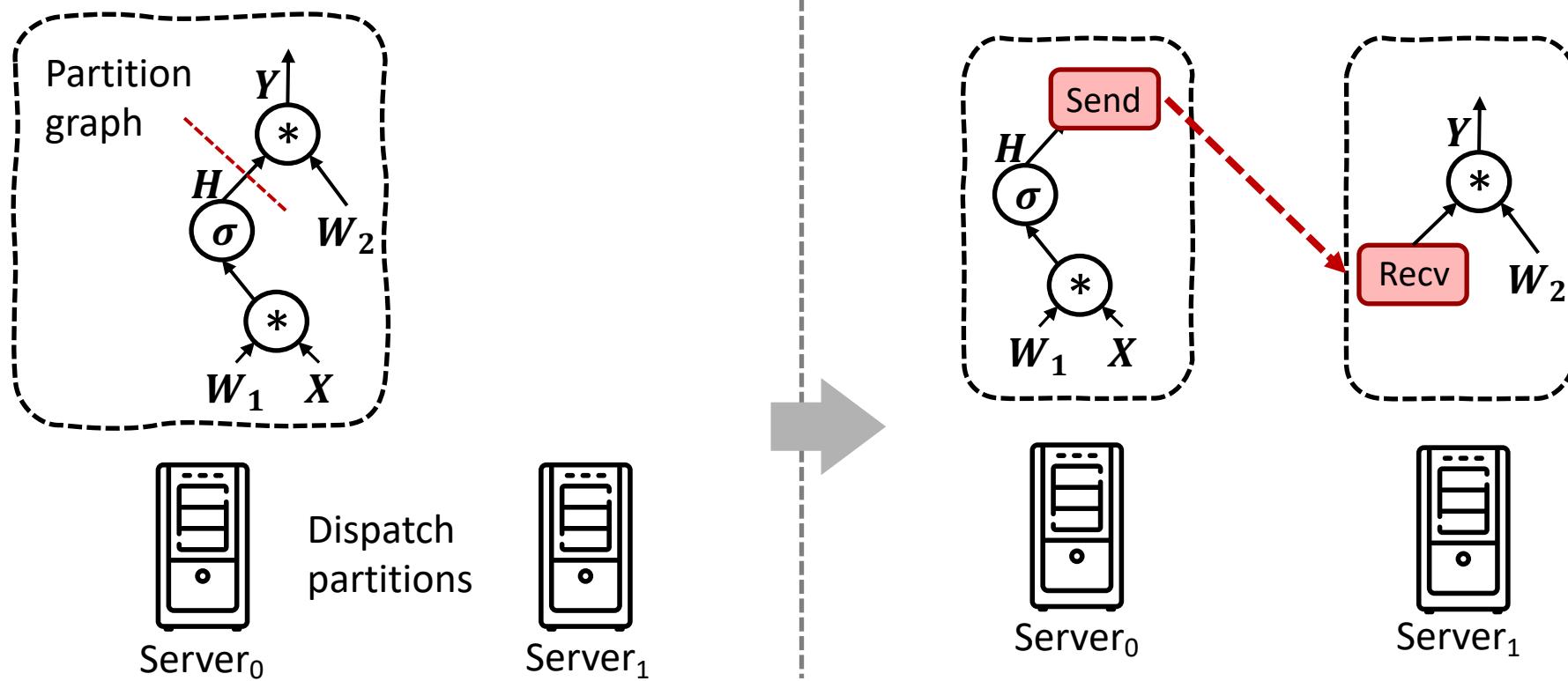
## 2. 具体内容

- **AI集群服务器架构**：参数服务器模式 – 同步与异步并行 - 环同步算法
- **AI集群软硬件通信**：通信软硬件实现 - 通信实现方式
- **分布式通信原语**：通信原语
- **框架分布式功能**：并行处理硬件架构 – AI框架中的分布式训练

- **大模型算法**：挑战 – 算法结构 – SOTA大模型
- **分布式并行**：数据并行 – 张量并行 – 自动并行 – 多维混合并行

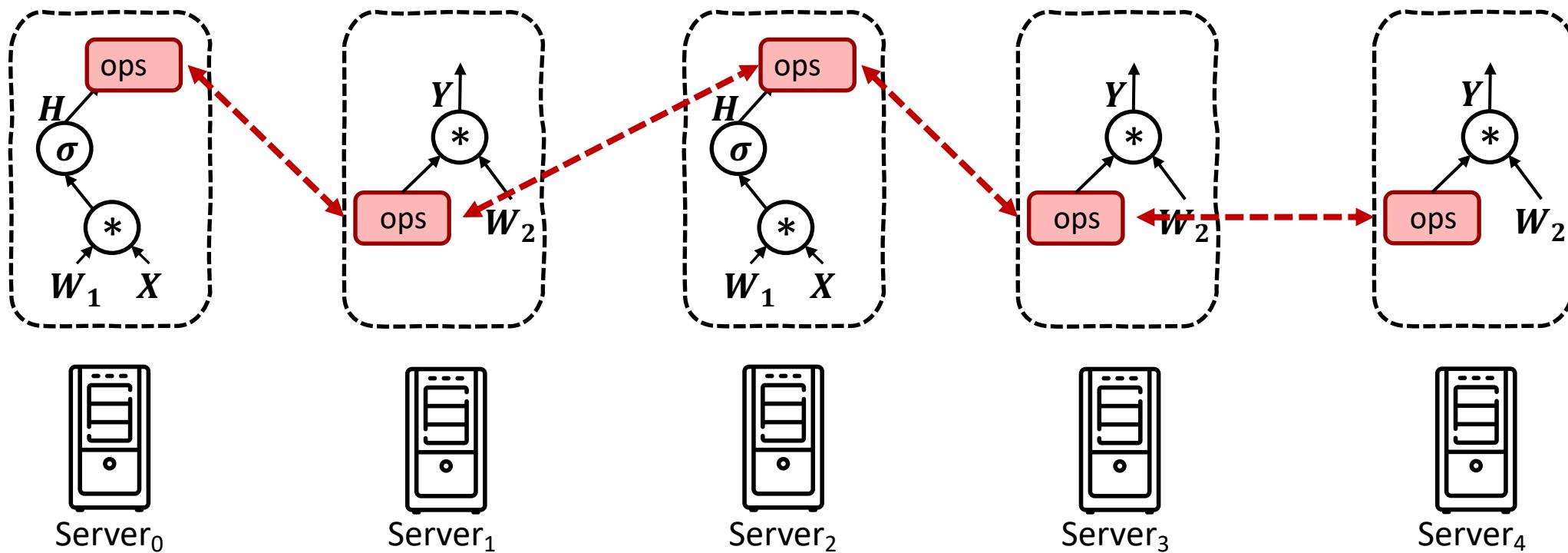
# 分布式训练系统

- 计算图跨节点切分



# 分布式训练系统

- 计算图跨节点同步数据（权重weights）信息



# 分布式训练系统 ( II )



- 示例：分布式MNIST

```
""" Gradient averaging. """
def average_gradients(model):
    size = float(dist.get_world_size())
    for param in model.parameters():
        dist.all_reduce(param.grad.data, op=dist.reduce_op.SUM)
        param.grad.data /= size
```

```
""" Distributed Synchronous SGD Example """
def run(rank, size):
    torch.manual_seed(1234)
    train_set, bsz = partition_dataset()
    model = Net()
    optimizer = optim.SGD(model.parameters(),
                          lr=0.01, momentum=0.5)

    num_batches = ceil(len(train_set.dataset) / float(bsz))
    for epoch in range(10):
        epoch_loss = 0.0
        for data, target in train_set:
            optimizer.zero_grad()
            output = model(data)
            loss = F.nll_loss(output, target)
            epoch_loss += loss.item()
            loss.backward()
            average_gradients(model)
            optimizer.step()
        print('Rank ', dist.get_rank(), ', epoch ',
              epoch, ': ', epoch_loss / num_batches)
```

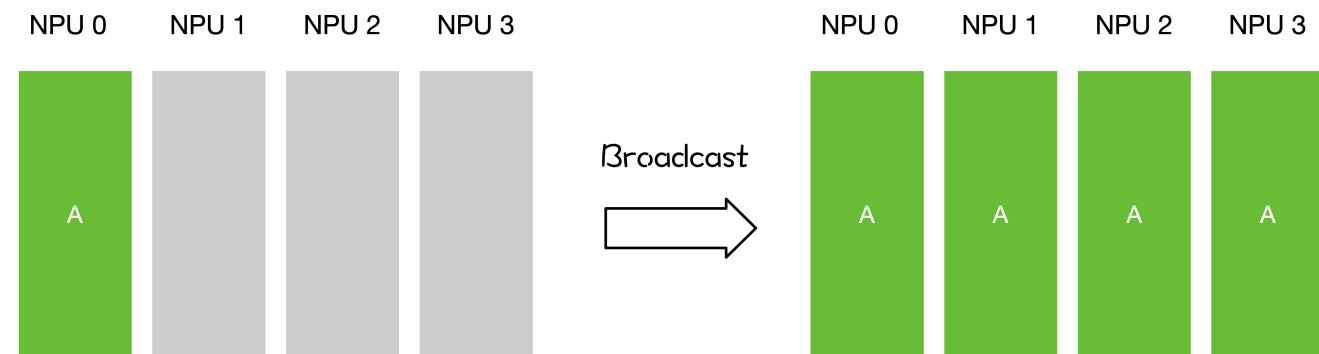
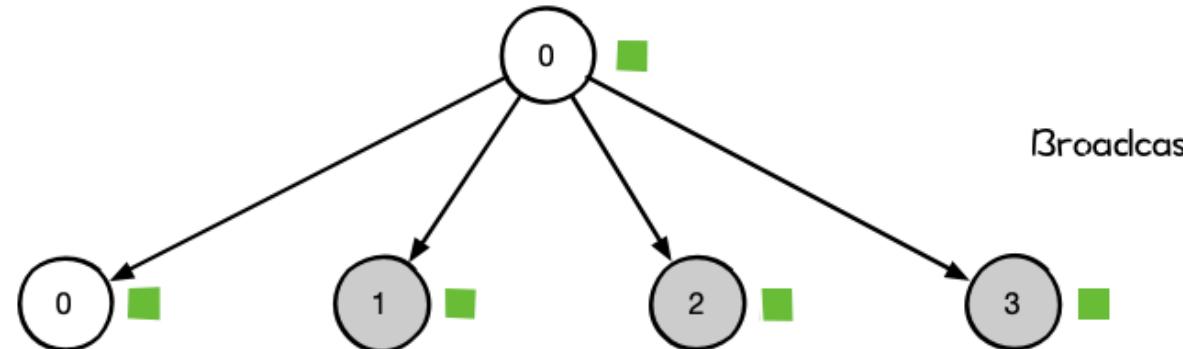
# 集合式通信方式

- 一对多 : Scatter / Broadcast
- 多对一 : Gather / Reduce
- 多对多 : All-Reduce / All-Gather

# 集合式通信方式 (I) : 一对多 Broadcast

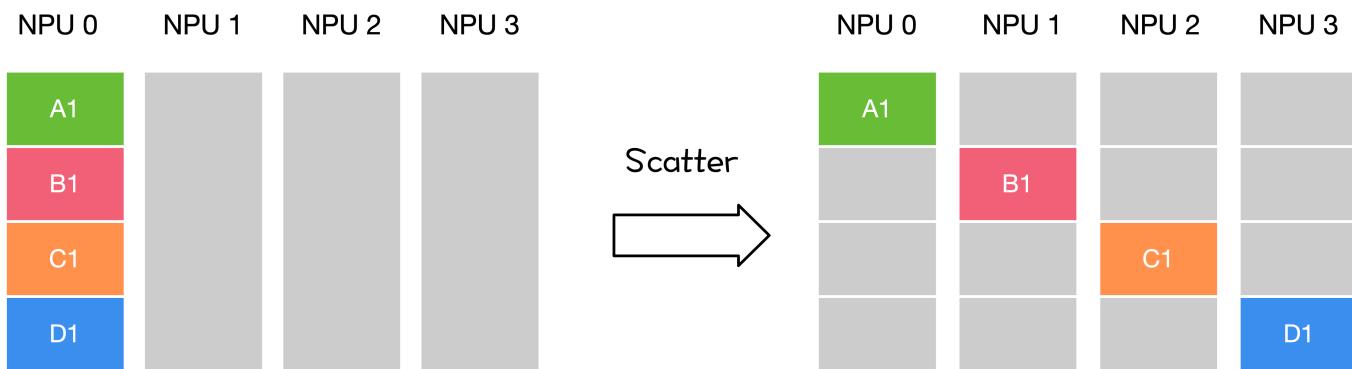
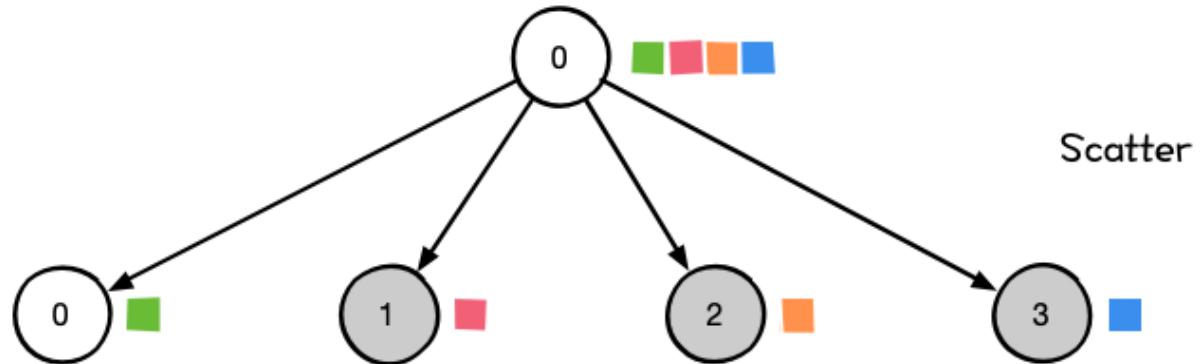
某个节点想把自身的数据发送到集群中的其他节点，那么就可以使用广播Broadcast的操作。

分布式机器学习中常用于网络参数的初始化。



# 集合式通信方式 (I) : 一对多 Scatter

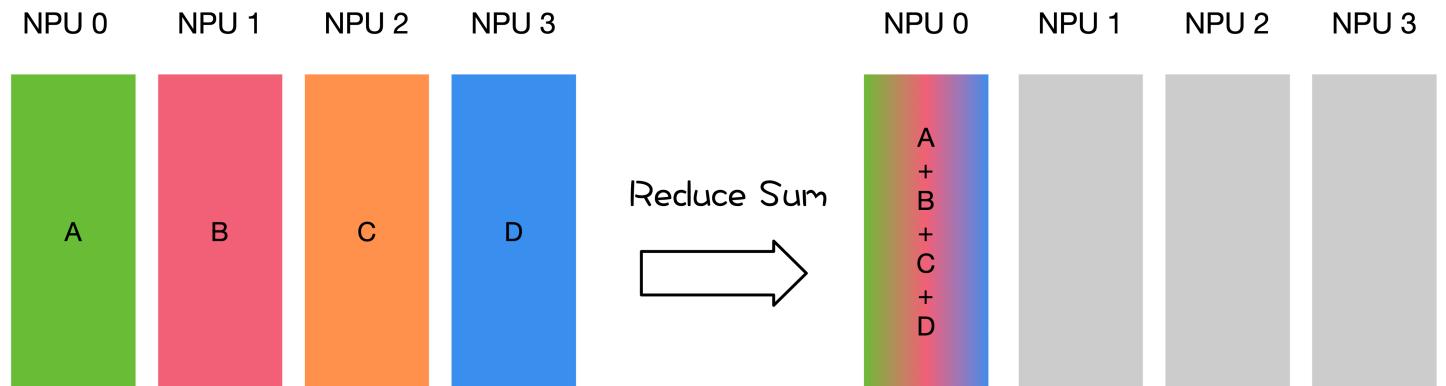
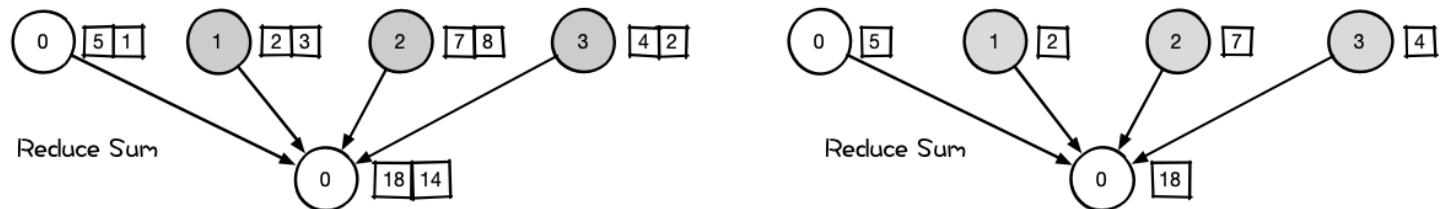
将主节点的数据进行划分并散布至其他指定的节点。



## 集合式通信方式 (II) : 多对一 Reduce

Reduce 称为规约运算，是一系列简单运算操作的统称。

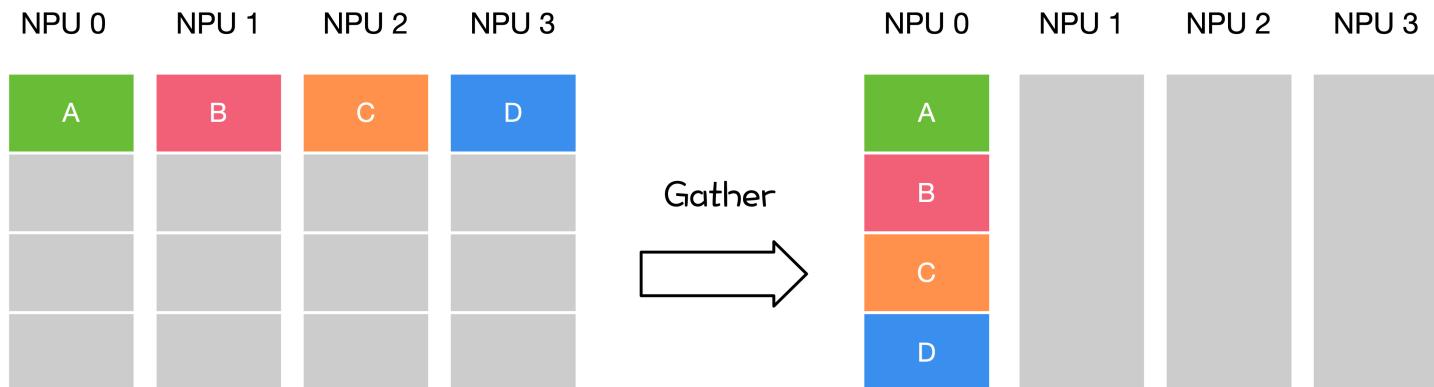
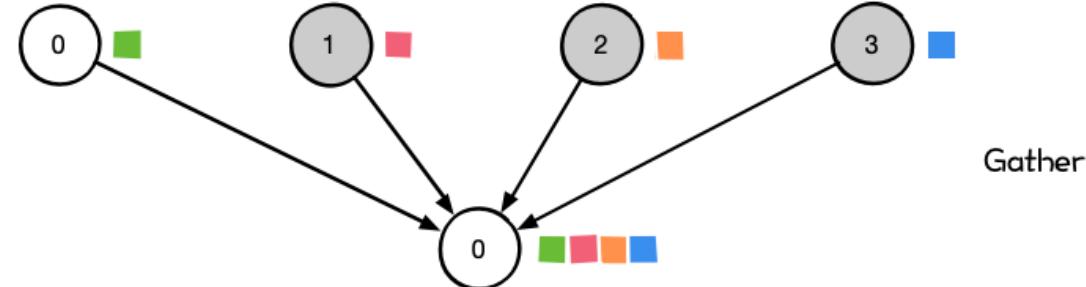
细分可以包括：SUM、MIN、MAX、PROD、LOR等类型的规约操作。



## 集合式通信方式 (II) : 多对一 Gather

将多个 Sender 上的数据收集到单个节点上，Gather 可以理解为反向的 Scatter。

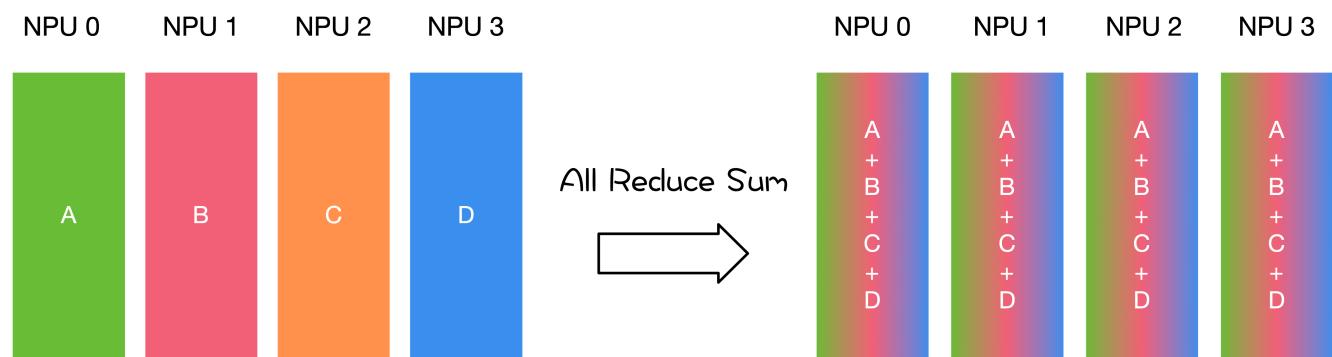
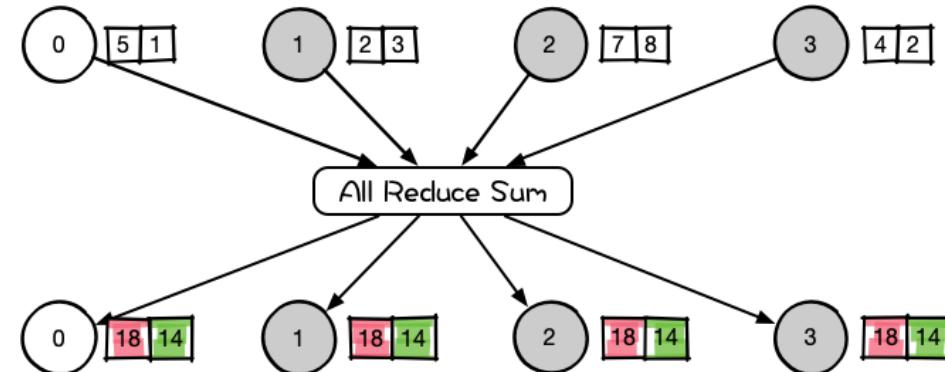
对很多并行算法很有用，比如并行的排序和搜索。



## 集合式通信方式 (III) : 多对多 All Reduce

All Reduce则是在所有的节点上都应用同样的Reduce操作。

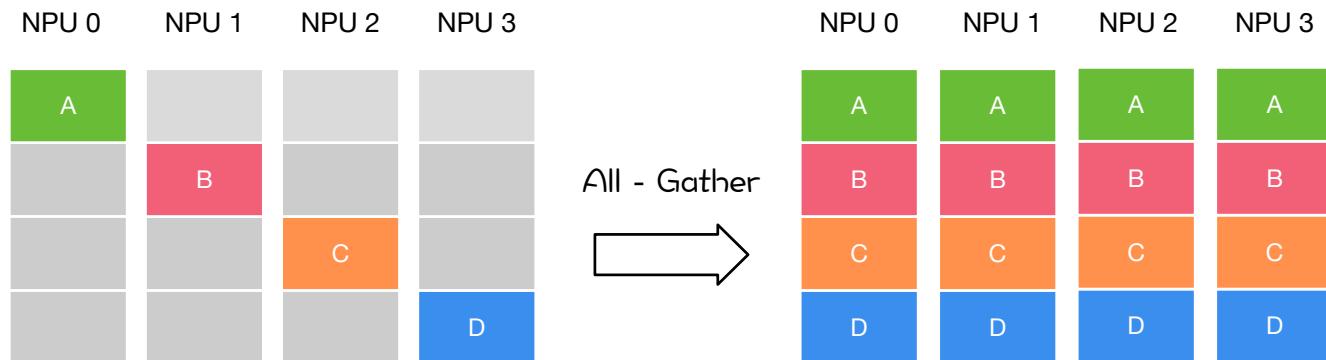
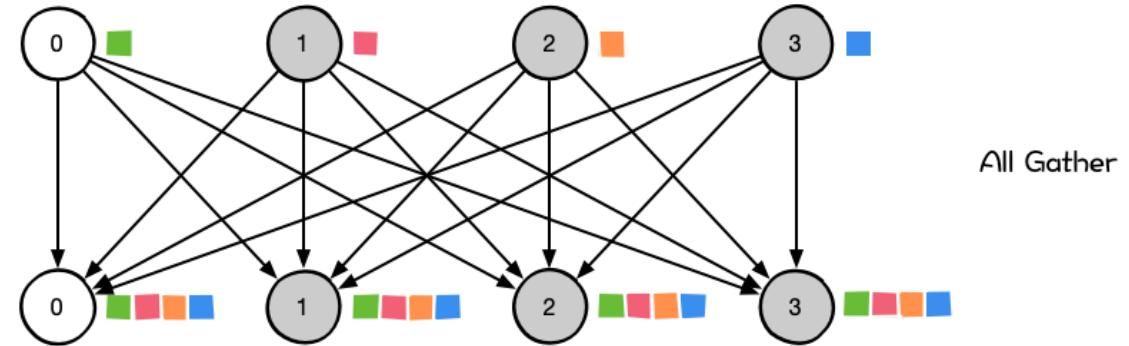
All Reduce操作可通过Reduce + Broadcast 或者 Reduce-Scatter + All-Gather 操作完成。



## 集合式通信方式 (III) : 多对多 All-Gather

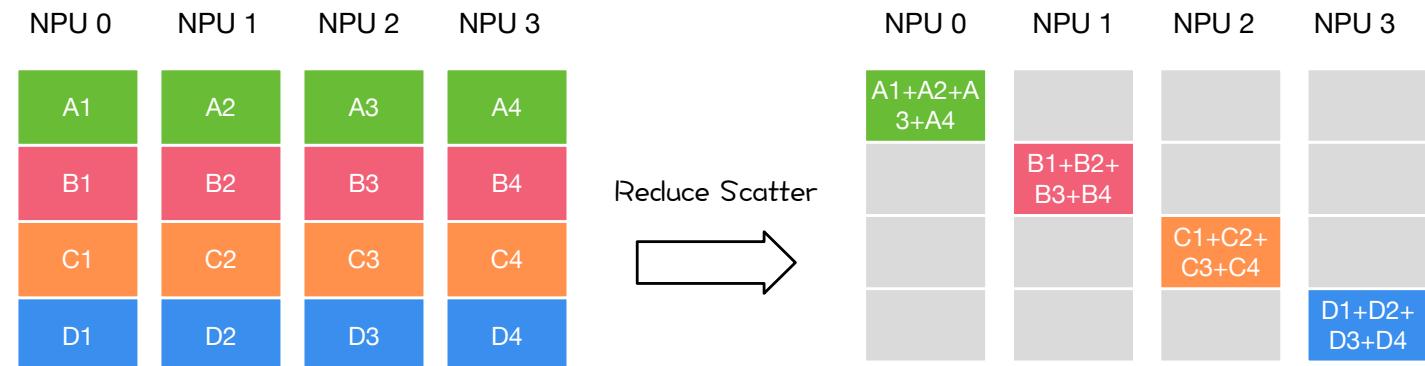
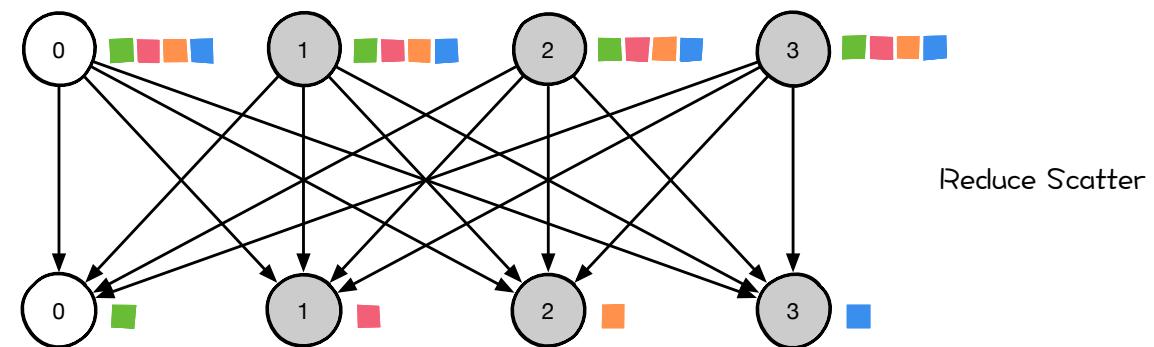
All Gather会收集所有数据到所有节点上。All Gather = Gather + Broadcast。

发送多个元素到多个节点很有用，即在多对多通信模式的场景。



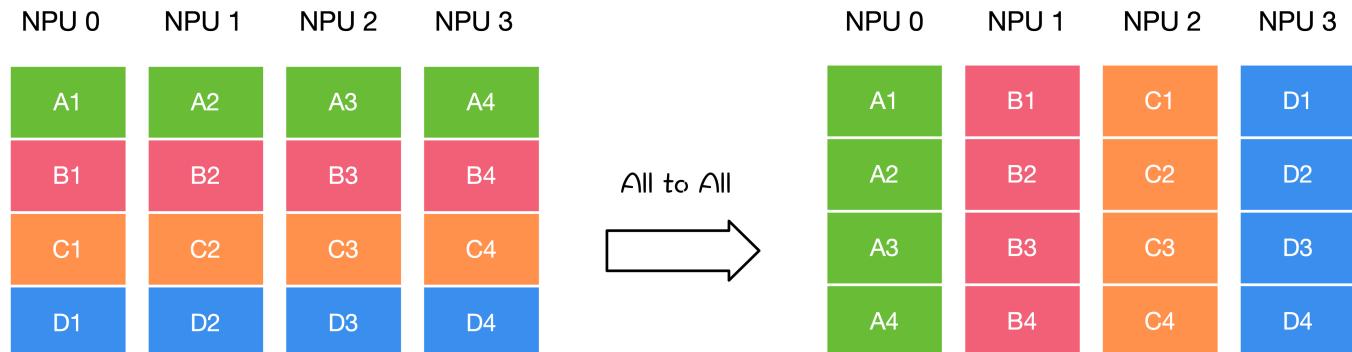
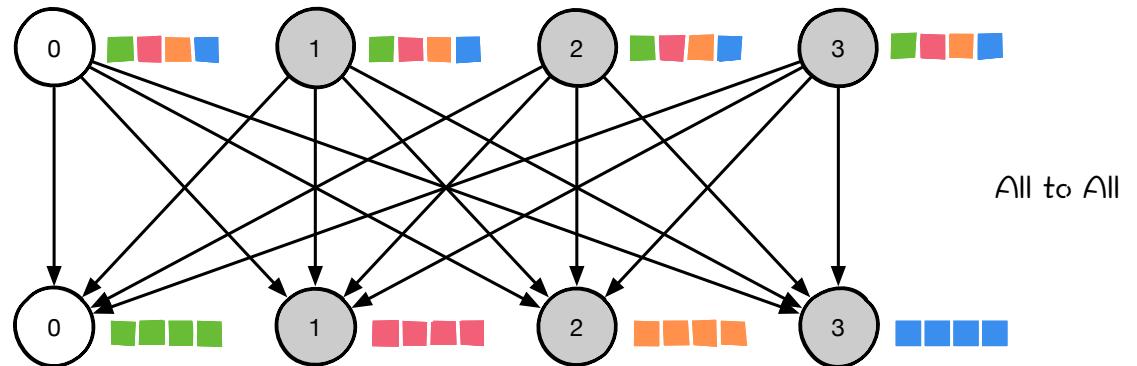
# 集合式通信方式 (III) : 多对多 Reduce Scatter

Reduce Scatter操作会将个节点的输入先进行求和，然后在第0维度按卡数切分，将数据分发到对应的卡上。

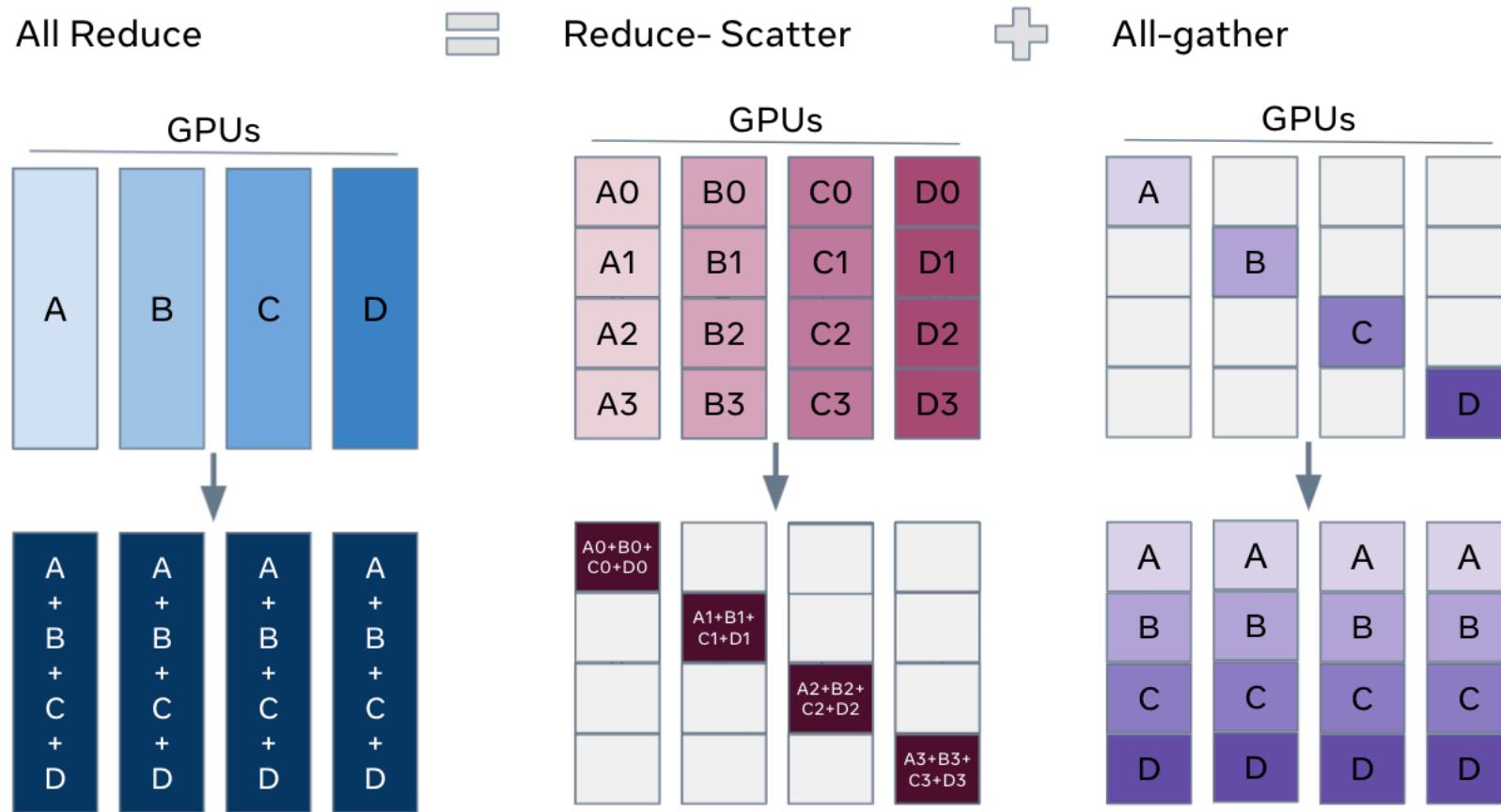


## 集合式通信方式 (III) : 多对多 All to All

将节点i的发送缓冲区中的第j块数据发送给节点j，节点j将接收到的来自节点i的数据块放在自身接收缓冲区的第i块位置。



# 集合式通信方式 (III) : 操作分解



# Summary

1. 了解集合式通信的3种不同方式
2. 了解一对多Scatter/Broadcast，多对一Gather/Reduce，多对多的具体方式
3. 了解多对多可以由一对多和多对一的方式组合， $\text{all-Reduce} = \text{Reduce Scatter} + \text{All gather}$



BUILDING A BETTER CONNECTED WORLD

THANK YOU

Copyright©2014 Huawei Technologies Co., Ltd. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.