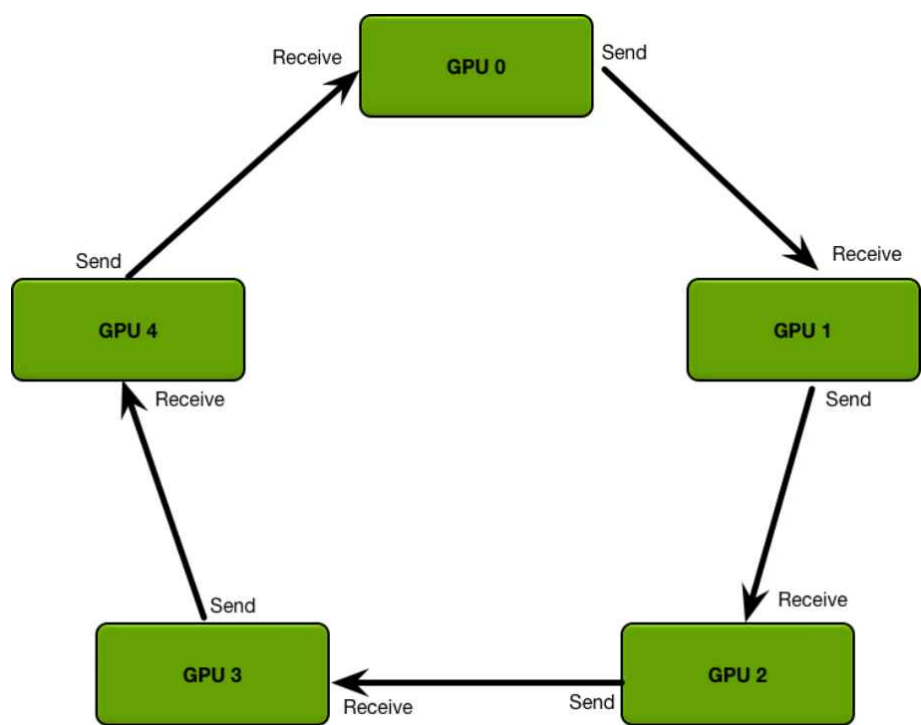


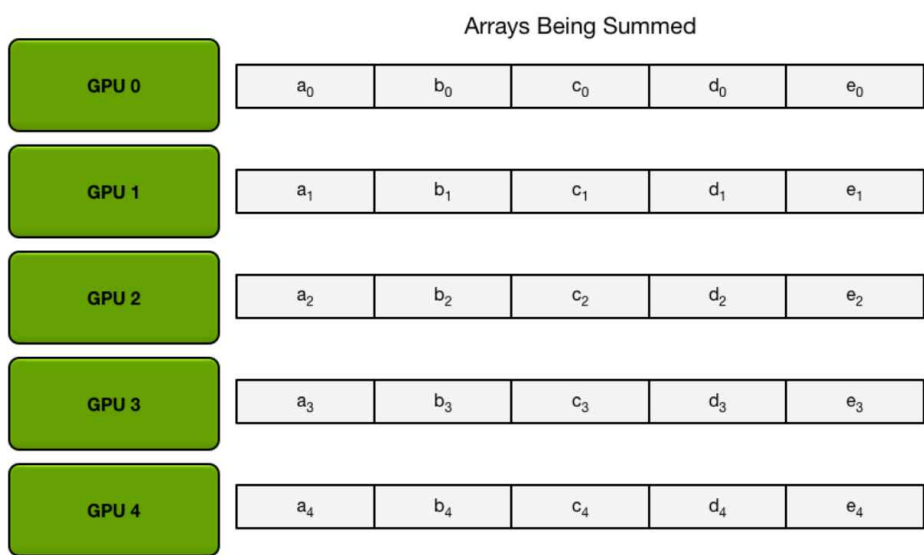
大模型 DeepSpeed 配置详解

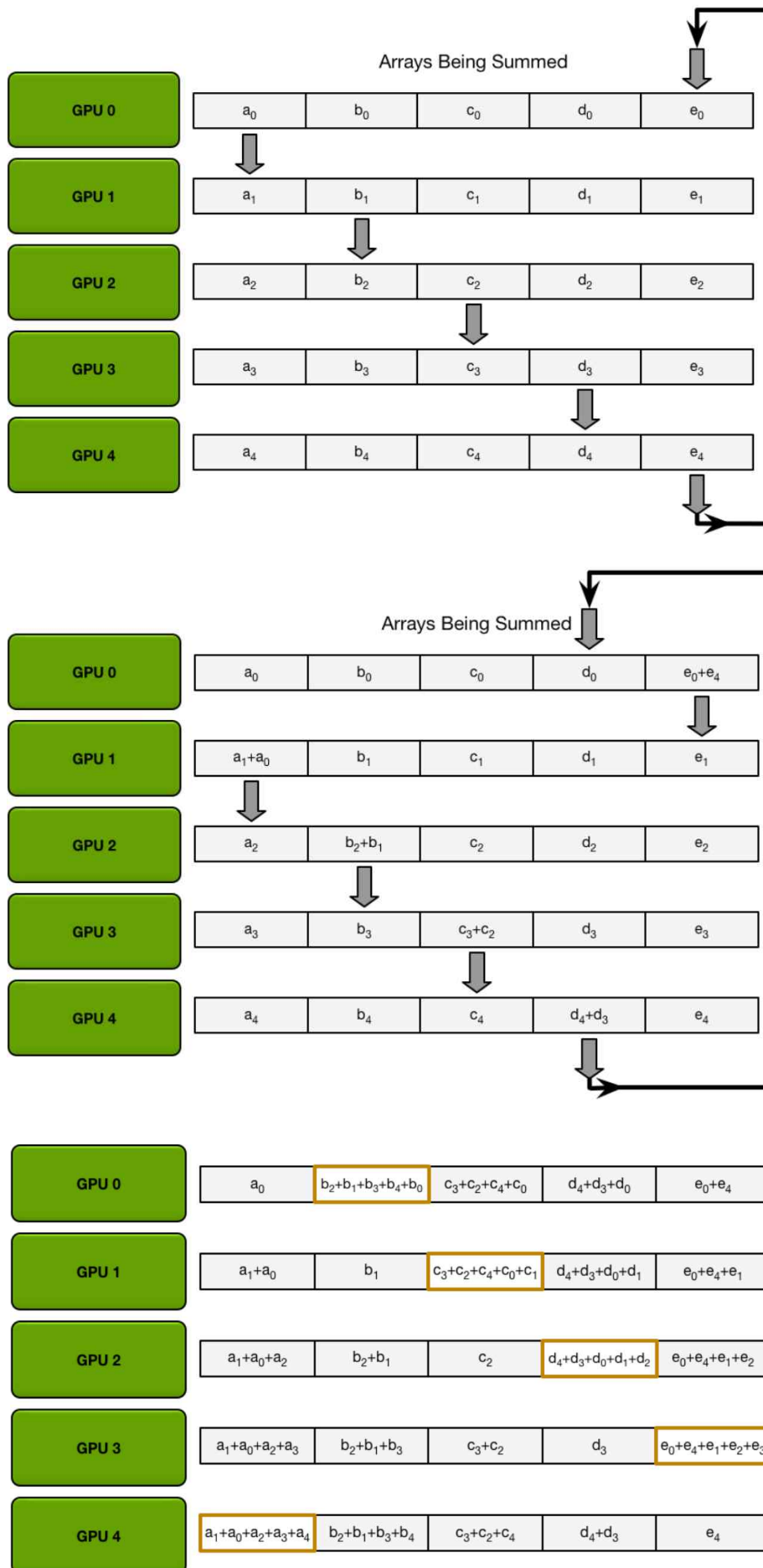
1. Stage2

【Ring AllReduce】

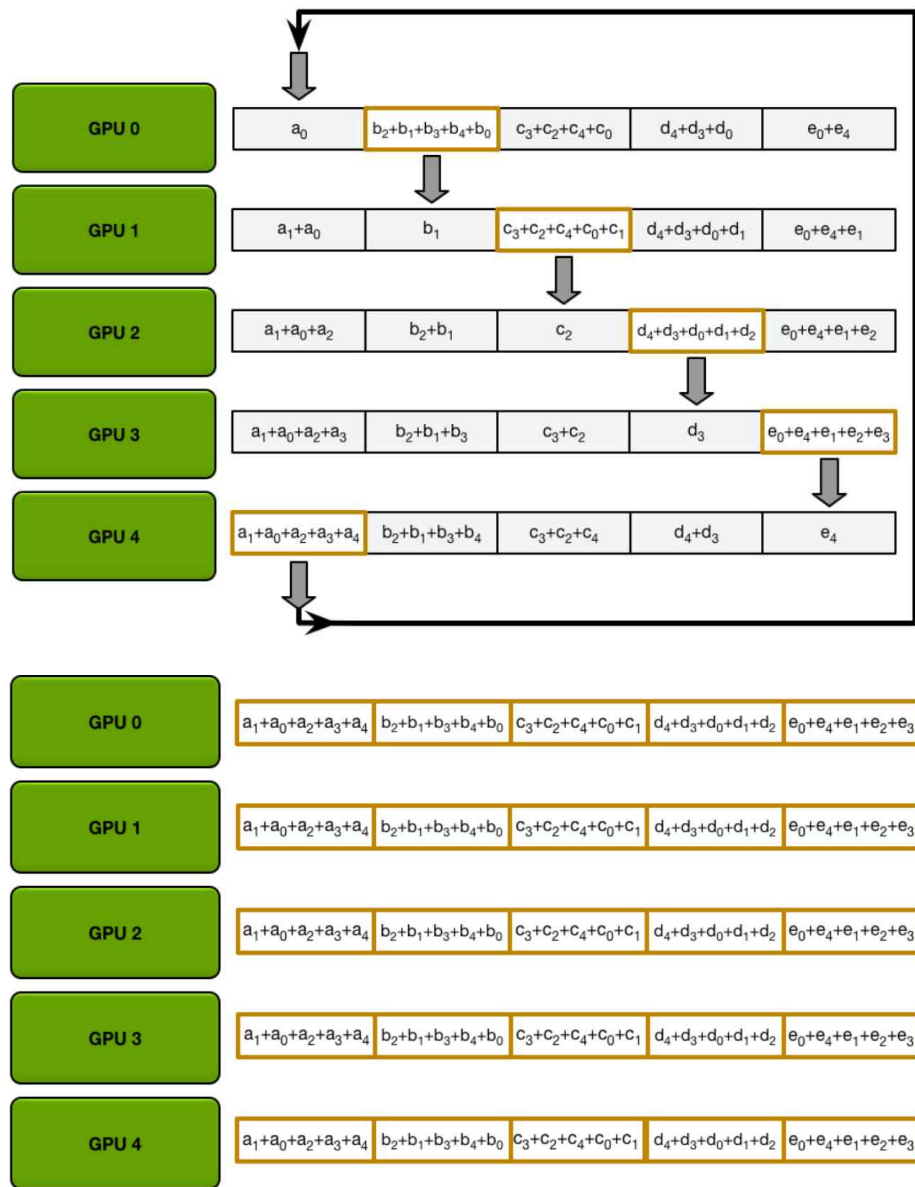


【Scatter-Reduce】





[The Allgather]



```

1 {
2     "optimizer": {
3         "type": "AdamW",
4         "params": {
5             "lr": "auto",
6             "weight_decay": "auto",
7             "torch_adam": true,
8             "adam_w_mode": true
9         }
10    },
11    "scheduler": {
12        "type": "WarmupDecayLR",
13        "params": {
14            "warmup_min_lr": "auto",
15            "warmup_max_lr": "auto",
16            "warmup_num_steps": "auto",
17            "total_num_steps": "auto"
18        }
19    }
20 }

```

```

19     },
20     "fp16": {
21         "enabled": true,
22         "loss_scale": 0,
23         "loss_scale_window": 1000,
24         "initial_scale_power": 16,
25         "hysteresis": 2,
26         "min_loss_scale": 1
27     },
28     "zero_optimization": {
29         "stage": 2,
30         "allgather_partitions": true,
31         "allgather_bucket_size": 2e8,
32         "reduce_scatter": true,
33         "reduce_bucket_size": "auto",
34         "overlap_comm": true,
35         "contiguous_gradients": true
36     },
37     "gradient_accumulation_steps": "auto",
38     "gradient_clipping": "auto",
39     "steps_per_print": 1000,
40     "train_batch_size": "auto",
41     "train_micro_batch_size_per_gpu": "auto",
42     "wall_clock_breakdown": false
43 }

```

```

1 {
2     "optimizer": {
3         # 优化器类型
4         "type": "AdamW",
5         # 优化器参数
6         "params": {
7             # 学习率
8             "lr": "auto",
9             # 学习率衰减
10            "weight_decay": "auto",
11            "torch_adam": true,
12            "adam_w_mode": true
13        }
14    },
15    # LR scheduler
16    "scheduler": {
17        "type": "WarmupDecayLR",
18        "params": {
19            "warmup_min_lr": "auto",

```

```

20         "warmup_max_lr": "auto",
21         "warmup_num_steps": "auto",
22         "total_num_steps": "auto"
23     }
24 },
25 # 使用利用 NVIDIA 的 Apex 包的 mixed precision/FP16 训练的配置
26 # 该模式类似于 AMP 的 O2 模式
27 # O0: 纯FP32训练, 可以作为accuracy的baseline。
28 # O1: 混合精度训练 (推荐使用), 根据黑白名单自动决定使用FP16 (GEMM, 卷积) 还是FP32
    (Softmax) 进行计算。
29 # O2: “几乎FP16”混合精度训练, 不存在黑白名单, 除了Batch norm, 几乎都是用FP16计算。
30 # O3: 纯FP16训练, 很不稳定, 但是可以作为speed的baseline;
31 "fp16": {
32     # 是否启用fp16训练
33     "enabled": true,
34     # fp16训练的损失缩放值
35     # 反向计算起始, Loss首先乘以Loss Scale值进行扩大, 避免反向梯度过小而产生下溢
36     # 这样在反向传播过程中所有值都扩大了相同倍数
37     # 设置为0.0时为 dynamic loss scaling
38     # 设置为其他值将会使用 static fixed loss scaling
39     "loss_scale": 0,
40     # 升高/降低动态损失比例值的窗口
41     "loss_scale_window": 1000,
42     # 初始 dynamic loss scale 的幂参数 (2^initial_scale_power)
43     "initial_scale_power": 16,
44     # 表示动态损耗缩放中的延迟偏移
45     "hysteresis": 2,
46     # 最小的 dynamic loss scale 数值
47     "min_loss_scale": 1
48 },
49 # 启用和配置 ZERO 优化, 与 FP16/BF16/FP32 和 Adam 优化器兼容。
50 "zero_optimization": {
51     # 选择 ZeRO 优化器的不同阶段。 Stage 0、1、2、3分别指禁用、优化器状态分区、优化
    器+梯度状态分区、优化器+梯度+参数分区。
52     "stage": 2,
53     # 是否使用 reduce-scatter 算法来进行梯度聚合。
54     # Reduce-scatter 是一种并行计算中常用的算法, 可以将数据均匀地分发给各个节点, 并
    同时进行聚合操作, 从而实现高效的梯度聚合。这种算法可以减少全局通信的开销, 并提高训练的性
    能。
55     "reduce_scatter": true,
56     # 一次 reduce 的数据量
57     "reduce_bucket_size": "auto",
58     # 使用 allgather 在每个步骤结束时从所有 GPU 收集更新的参数
59     "allgather_partitions": true,
60     # 一次 allgather 的数据量
61     "allgather_bucket_size": "auto",
62     # 尝试将梯度的更新与反向传播重叠进行

```

```

63     "overlap_comm": true,
64     # 在生成梯度时将梯度复制到连续缓冲区。实现高效访问数据，避免反向传播期间的内存碎片。
65     # 通过设置 contiguous_gradients 参数为 true, DeepSpeed 可以要求梯度在内存中是连续存储的。这样做可以确保梯度聚合操作能够高效地进行，减少了数据拷贝和内存访问的开销。
66     # 相反，当 contiguous_gradients 设置为 false 时, DeepSpeed 不要求梯度在内存中是连续存储的。这意味着梯度聚合操作可能需要进行额外的数据拷贝和内存访问，导致性能下降。
67     "contiguous_gradients": true
68 },
69 # 梯度累计步数
70 # 在平均和使用梯度之前对梯度进行累计的步数，
71 # ① 会导致步骤之间的梯度通信频率降低 ② 能够使用每个 GPU 进行更大的批量训练
72 "gradient_accumulation_steps": "auto",
73 # 梯度裁剪，避免梯度爆炸，对超出预设范围的梯度值，进行梯度裁剪
74 "gradient_clipping": "auto",
75 # 信息打印间隔
76 "steps_per_print": 1000,
77 # 训练 batch_size
78 # train_batch_size must be equal to train_micro_batch_size_per_gpu *
  gradient_accumulation * number of GPUs, 3个参数指定2/3即可
79 "train_batch_size": "auto",
80 # 每个 GPU 的 micro_batch_size
81 "train_micro_batch_size_per_gpu": "auto",
82 # DeepSpeed 会在训练结束时输出一个时间统计报告。该报告显示了整个训练过程中的各个步骤的耗时情况，以及每个步骤中的细分时间。
83 # 这些步骤可能包括数据加载、前向传播、反向传播、优化器更新等等。通过查看时间统计报告，你可以了解每个步骤所花费的时间，并且可以根据需要进行性能优化和调整。
84 "wall_clock_breakdown": false
85 }

```

2. Stage3

```

1  {
2      "fp16": {
3          "enabled": true,
4          "auto_cast": false,
5          "loss_scale": 0,
6          "initial_scale_power": 32,
7          "loss_scale_window": 1000,
8          "hysteresis": 2,
9          "min_loss_scale": 1
10     },
11     "optimizer": {
12         "type": "AdamW",
13         "params": {

```

```

14         "lr": "auto",
15         "betas": [
16             0.9,
17             0.999
18         ],
19         "eps": 1e-8,
20         "weight_decay": "auto"
21     }
22 },
23 "zero_optimization": {
24     "stage": 3,
25     "overlap_comm": true,
26     "contiguous_gradients": true,
27     "stage3_prefetch_bucket_size": "auto",
28     "stage3_param_persistence_threshold": "auto",
29     "stage3_max_live_parameters": 0,
30     "stage3_max_reuse_distance": 0,
31     "stage3_gather_16bit_weights_on_model_save": true
32 },
33 "train_batch_size": "auto",
34 "train_micro_batch_size_per_gpu": "auto",
35 "wall_clock_breakdown": false
36 }

```

```

1 {
2     "fp16": {
3         "enabled": true,
4         "auto_cast": false,
5         "loss_scale": 0,
6         "initial_scale_power": 32,
7         "loss_scale_window": 1000,
8         "hysteresis": 2,
9         "min_loss_scale": 1
10    },
11    "optimizer": {
12        "type": "AdamW",
13        "params": {
14            "lr": "auto",
15            "betas": [
16                0.9,
17                0.999
18            ],
19            "eps": 1e-8,
20            "weight_decay": "auto"
21        }
22    }
23 }

```

```

22     },
23     "zero_optimization": {
24         "stage": 3,
25         # 启用将优化器状态卸载到 CPU 或 NVMe, 并将优化器计算卸载到 CPU。 这可以为更大的
模型或批量大小释放 GPU 内存。 适用于 ZeRO 阶段 1、2、3。
26         "offload_optimizer": {
27             "device": "cpu",
28             "pin_memory": true
29         },
30         # 启用和配置参数卸载到 CPU 或 NVMe, 仅适用于 ZeRO stage 3。 请注意, 如果未指定
或不支持“device”的值, 则会触发异常。
31         "offload_param": {
32             "device": "cpu",
33             "pin_memory": true,
34             "max_in_cpu": 1.4e10
35         },
36         # 尝试将梯度的更新与反向传播重叠进行
37         "overlap_comm": true,
38         # 在生成梯度时将其复制到连续的缓冲区, 避免向后传递期间出现内存碎片。
39         "contiguous_gradients": true,
40         # 用于预取参数的固定缓冲区的大小
41         # 较小的值使用较少的内存, 但可能会因通信而增加停顿。
42         "stage3_prefetch_bucket_size": "auto",
43         # 不对小于这个阈值的参数进行partition
44         # 较小的数值使用的内存更少, 但可能会极大地增加通讯负担
45         "stage3_param_persistence_threshold": "auto",
46         # 每个 GPU 驻留的最大参数量
47         # 较小的数值使用的内存更少, 但需要进行更多的通讯
48         "stage3_max_live_parameters": 0,
49         # 如果参数在这个阈值范围内会被复用, 不会被释放
50         # 较小的数值使用的内存更少, 但需要进行更多的通讯
51         "stage3_max_reuse_distance": 0,
52         # 在通过 save_16bit_model() 保存模型之前合并权重
53         # 由于权重在 GPU 之间进行 partition, 启用此选项时该函数会自动收集权重, 然后保
存 fp16 模型权重。
54         "stage3_gather_16bit_weights_on_model_save": true
55     },
56     "train_batch_size": "auto",
57     "train_micro_batch_size_per_gpu": "auto",
58     "wall_clock_breakdown": false
59 }

```