

新型总线技术研究 – Scale UP&OUT – UB

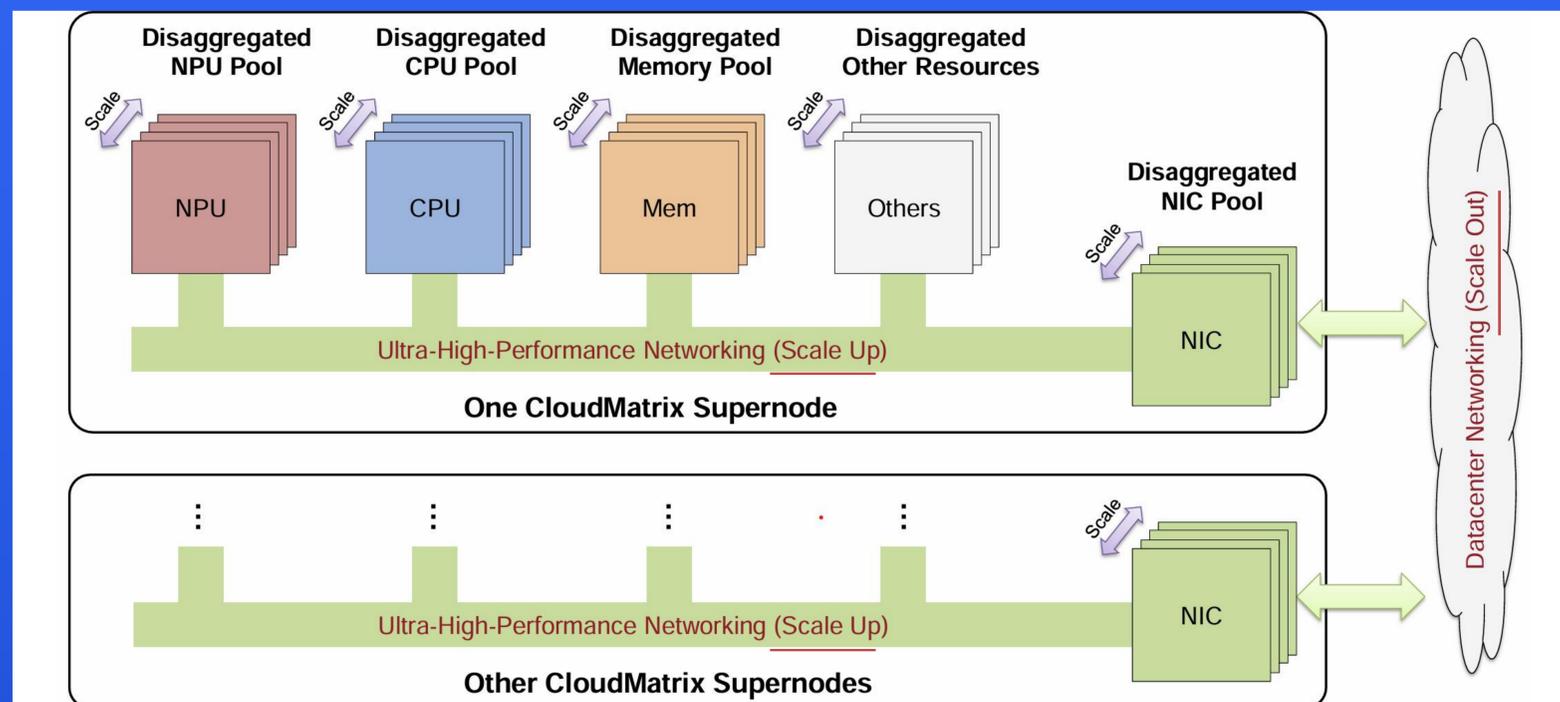
黎亮, openEuler Sig-Long committer, 2025-08-22

目录

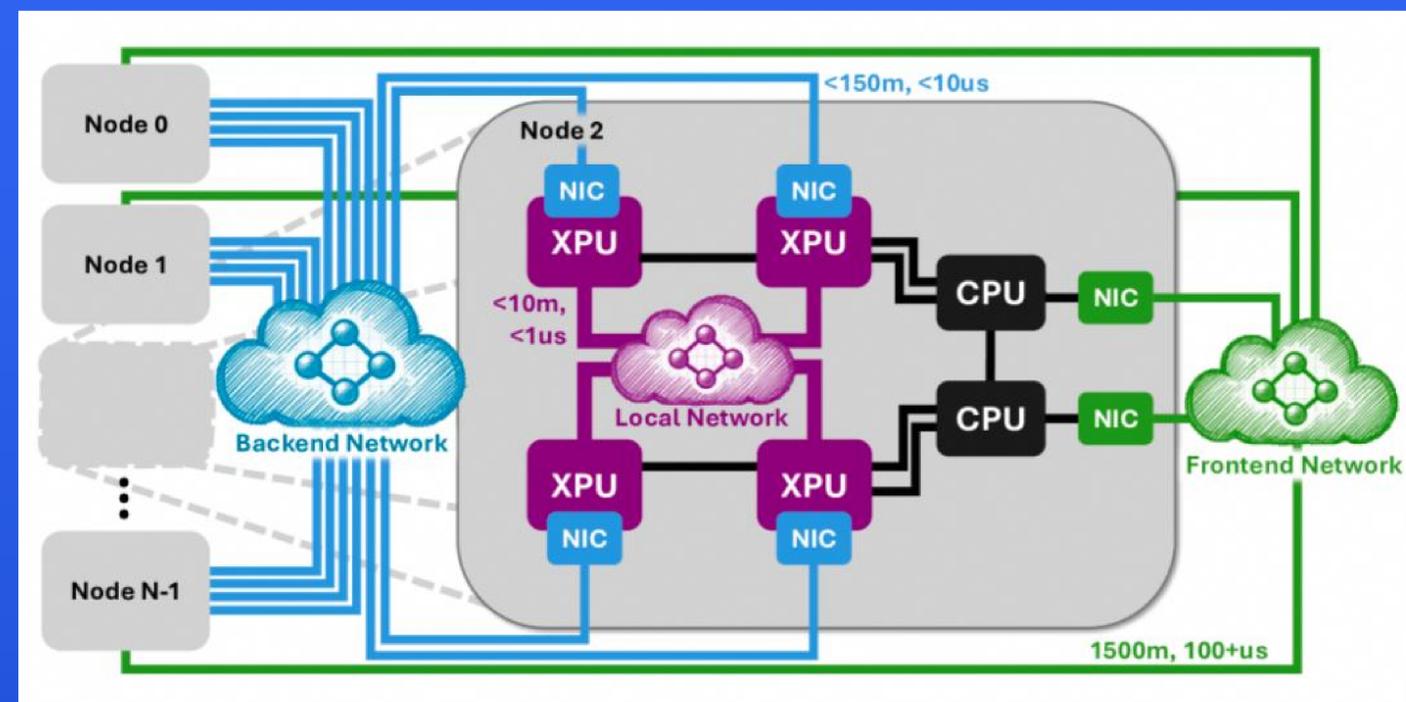
- 大规模异构计算系统的互联架构
- Scale UP 总线技术观察
- UB – Scale Up & Out 融合总线
- 内存语义与编程范式讨论
- 讨论&展望

大规模异构计算系统的互联架构

(大规模) 异构计算系统



<https://arxiv.org/pdf/2503.20377>



<https://arxiv.org/pdf/2506.12708>

Scale UP / OUT 边界区分的维度:

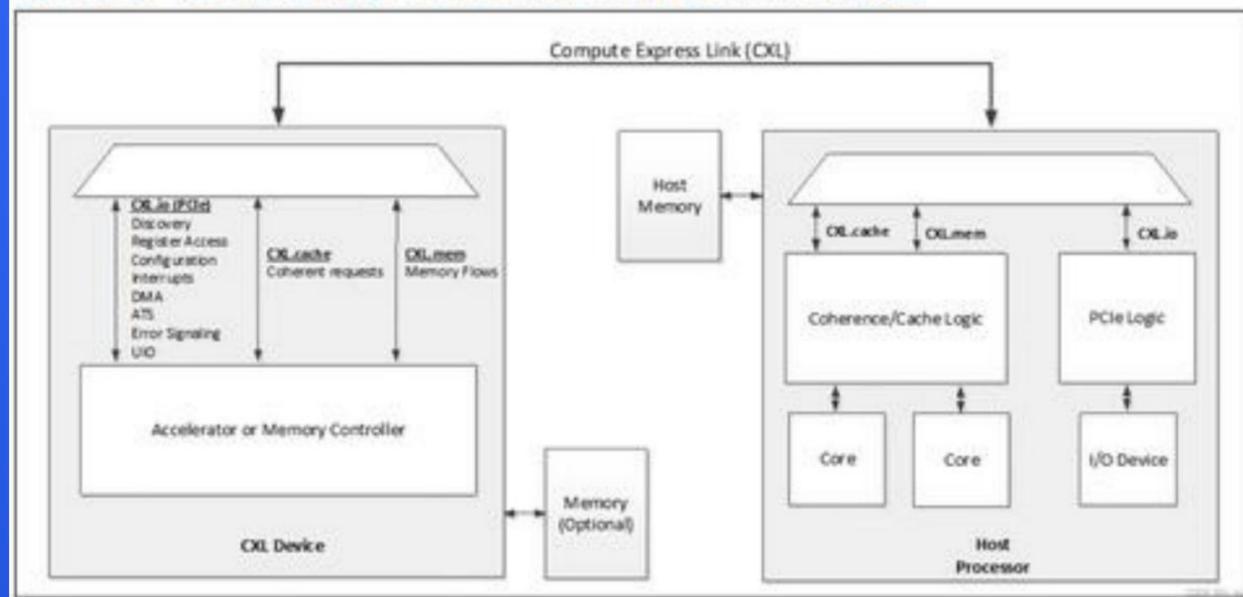
- 互联介质
- 物理距离
- 传输带宽
- 访问时延
- 编/寻址与数据移动方式
- 可靠性
- tbd

Topic: 当基于内存语义时, 边界还存在吗?

Scale UP 总线技术观察

CXL, NVLink, Scale-UP Ethernet, UB -- <1>

Figure 1-1. Conceptual Diagram of Device Attached to Processor via CXL



Fully Connected NVLink across 256 GPUs

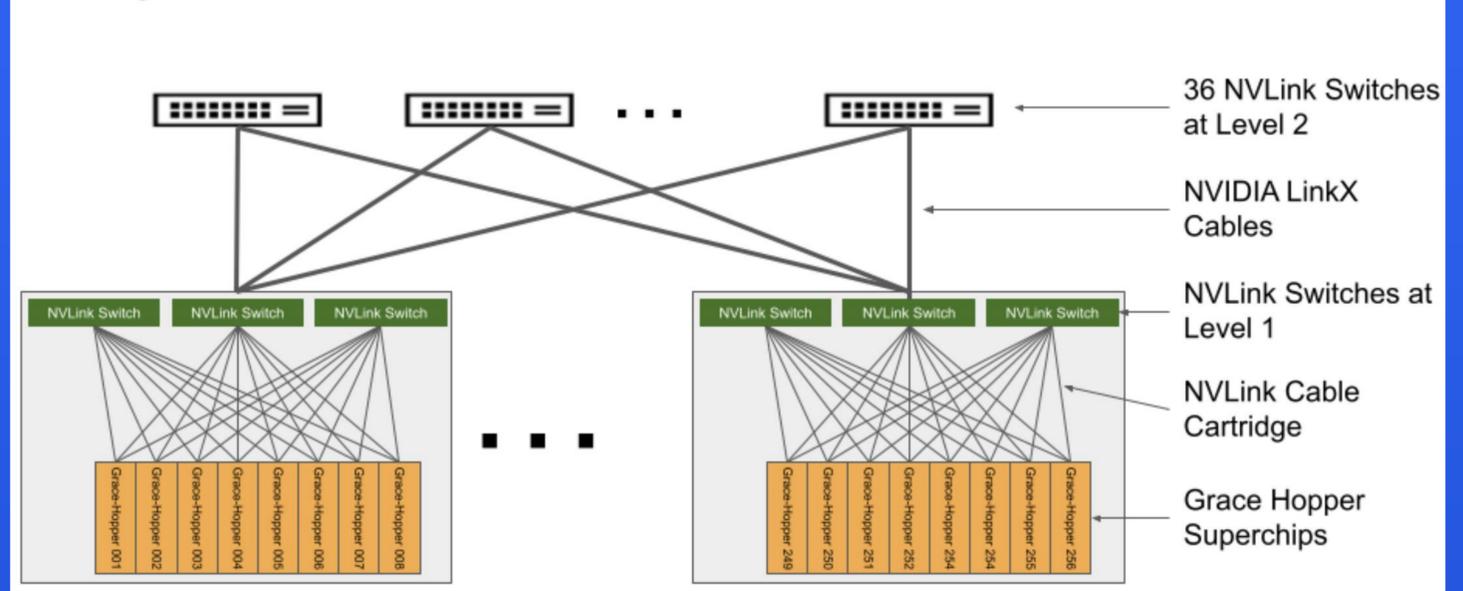
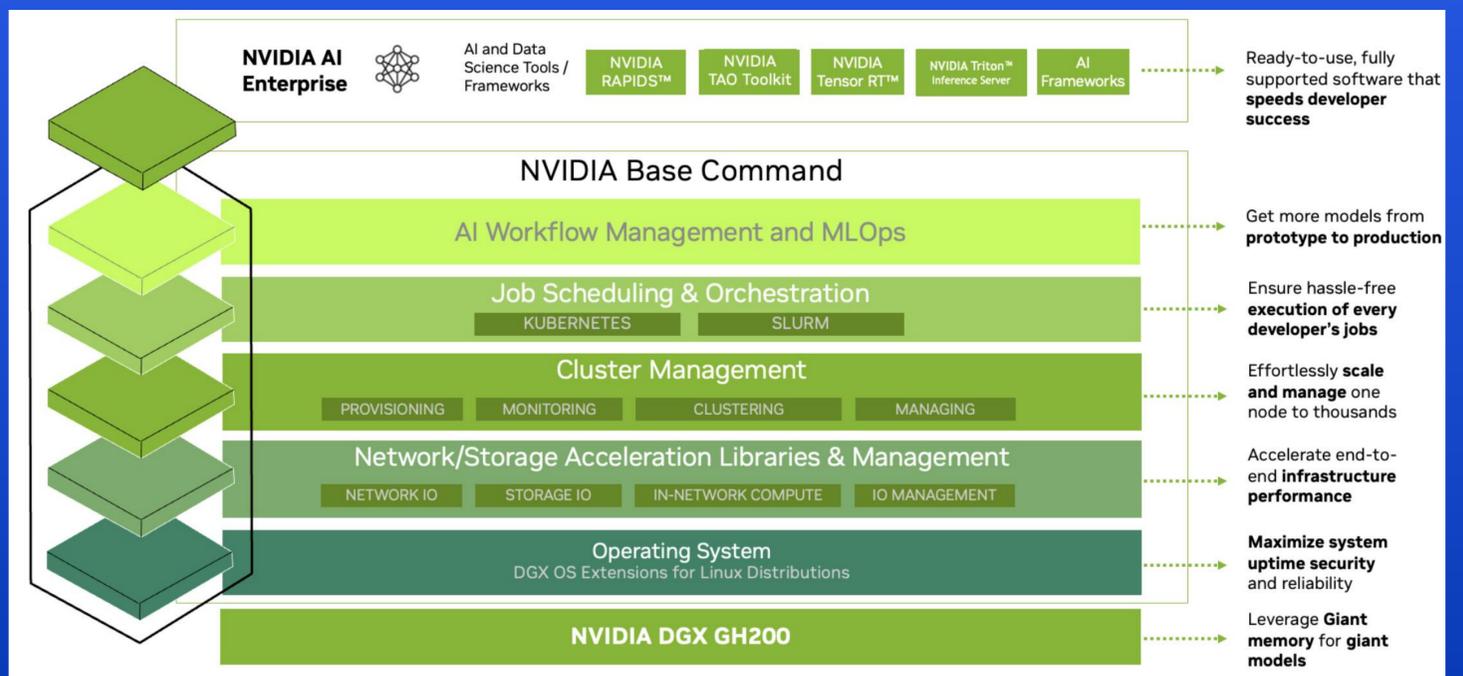
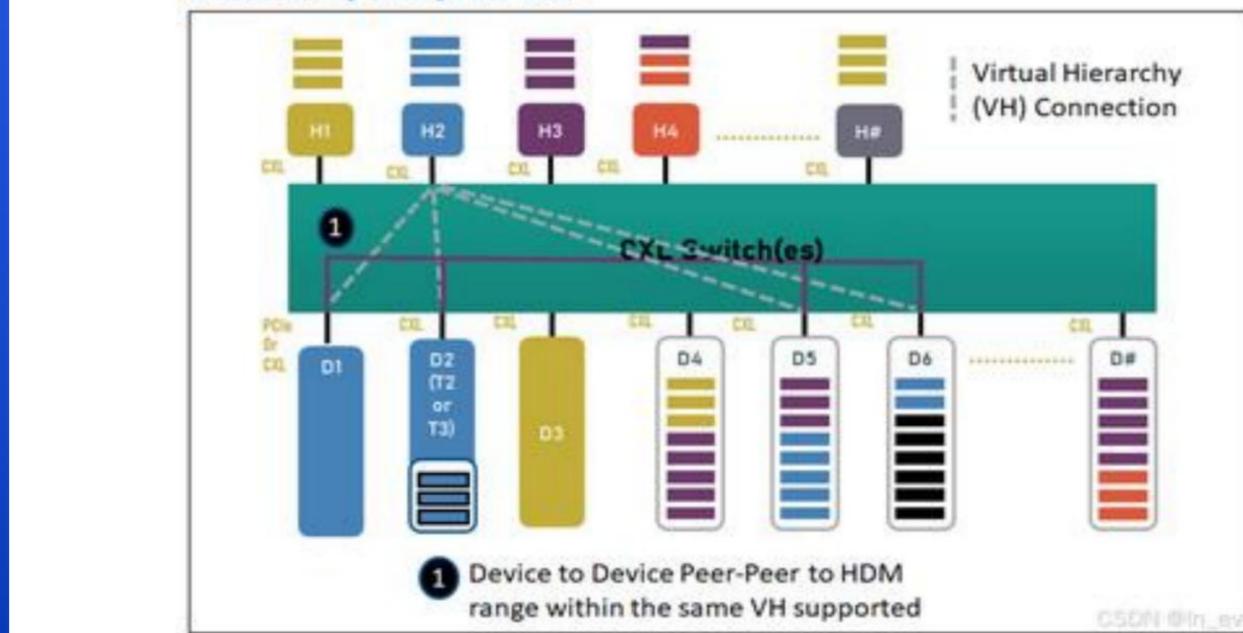
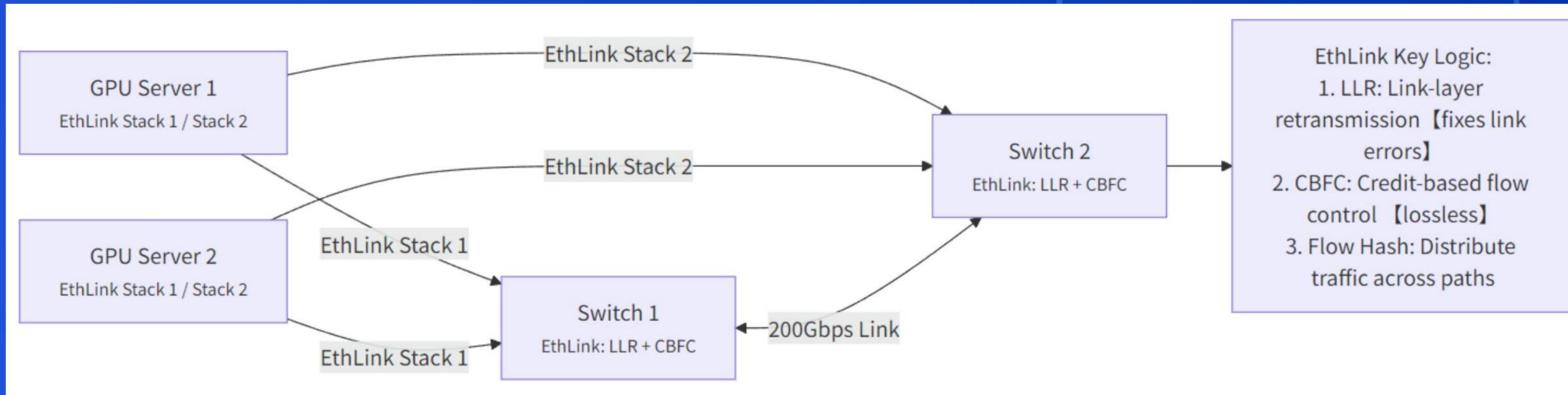


Figure 1-3. Direct Peer-to-Peer Access to an HDM Memory by PCIe/CXL Devices without Going through the Host



CXL, NVLink, Scale-UP Ethernet, UB —<2>

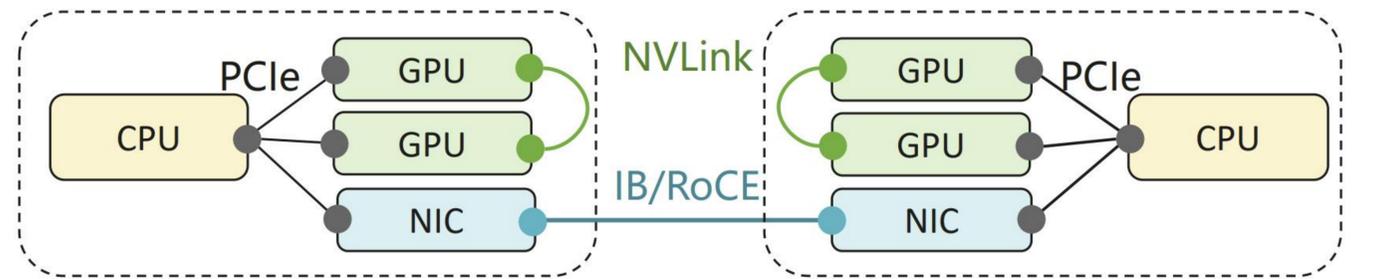


维度	EthLink	ETH-X	SUE
主导者	字节跳动 (企业级)	ODCC (行业联盟)	博通
节点规模	单域 1024+ (横向扩展)	256 卡全互联 (中规模)	10-100 (纵向扩展)
协议复杂度	双层优化 (保留传统层)	贴近标准以太网 (易兼容)	极简链路 + 事务层 (无 IP/TCP)
软件依赖	需上层处理乱序	统一编址简化软件	完全硬件化 (零软件)
硬件要求	标准交换机 + 多接口 NIC	支持 PAXI 的智能交换机	普通交换机 + 标准 PHY
典型延迟	亚微秒级 (100-500ns)	1-5 微秒	1-10 微秒

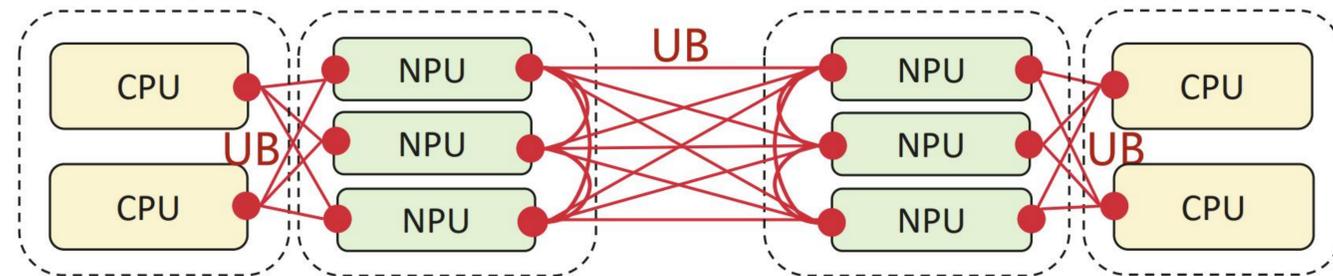
UB — to be announced

UB – Scale Up & Out 融合总线

UB – Scale Up & Out - Topo

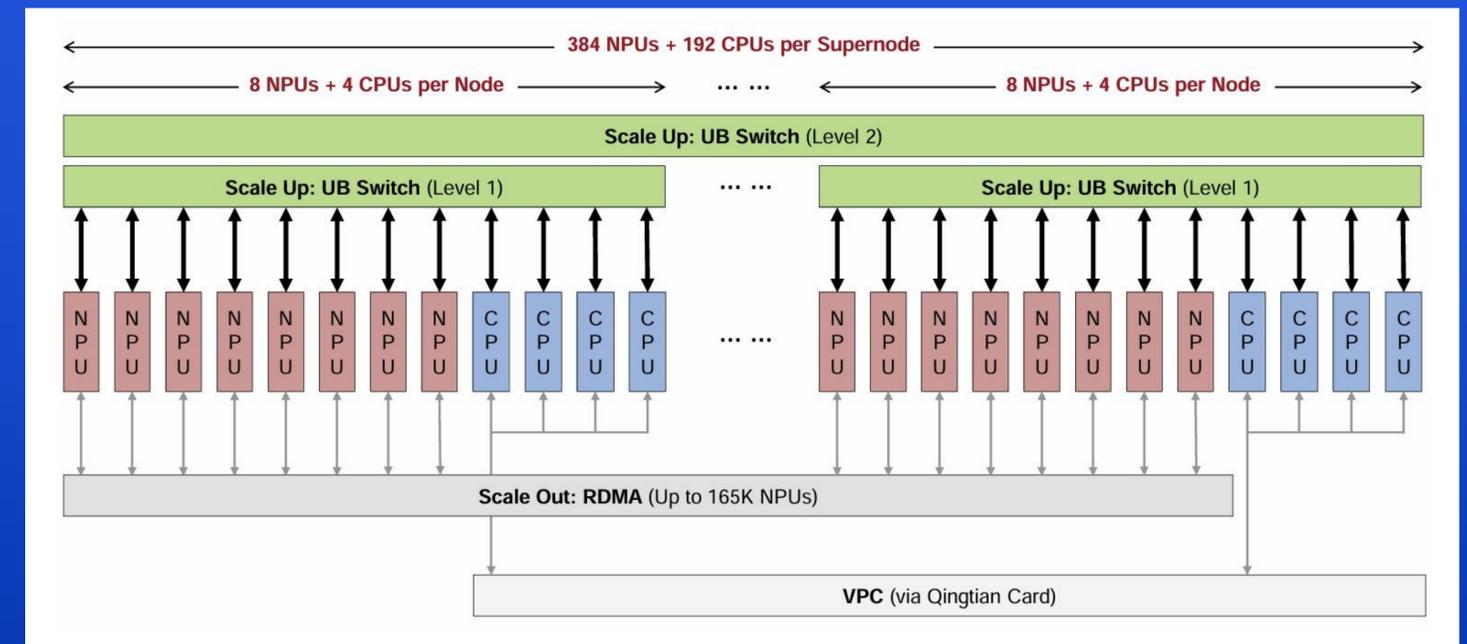
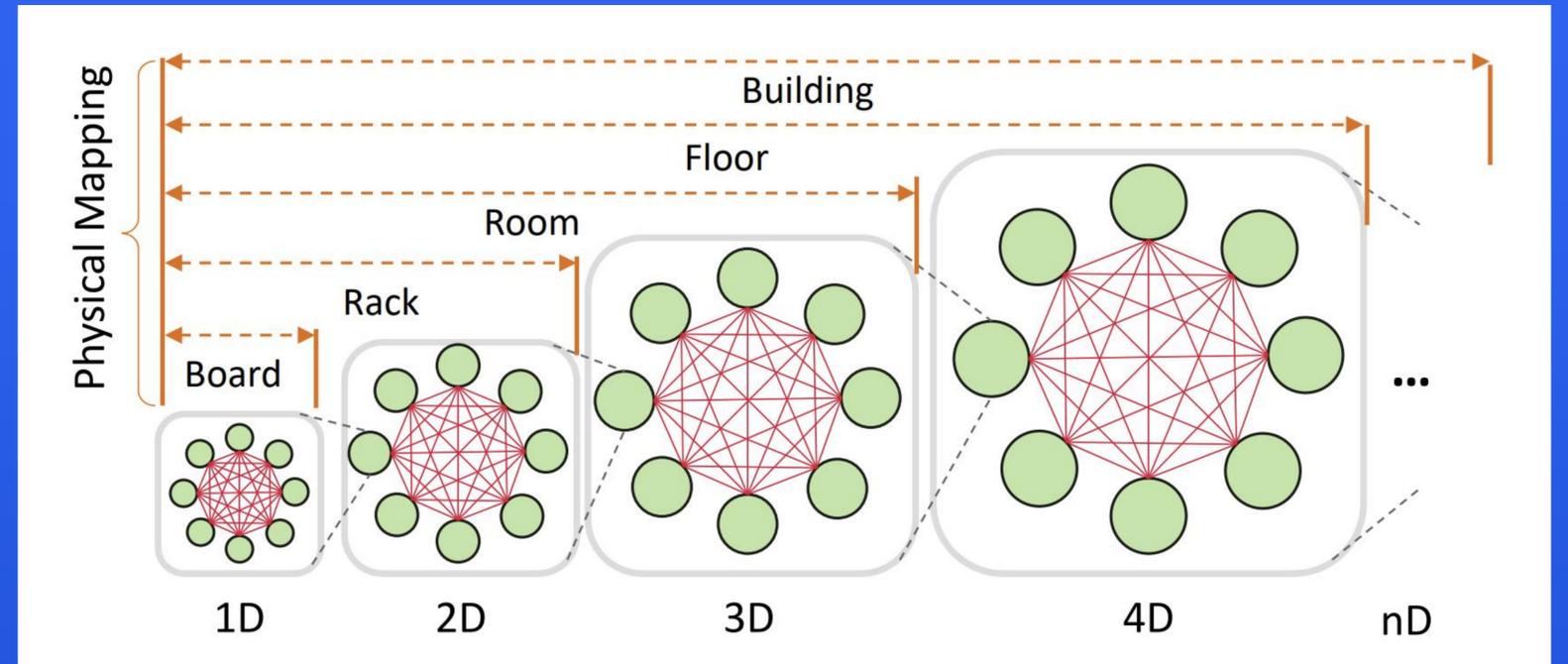


(a) Traditional architecture using hybrid interconnections

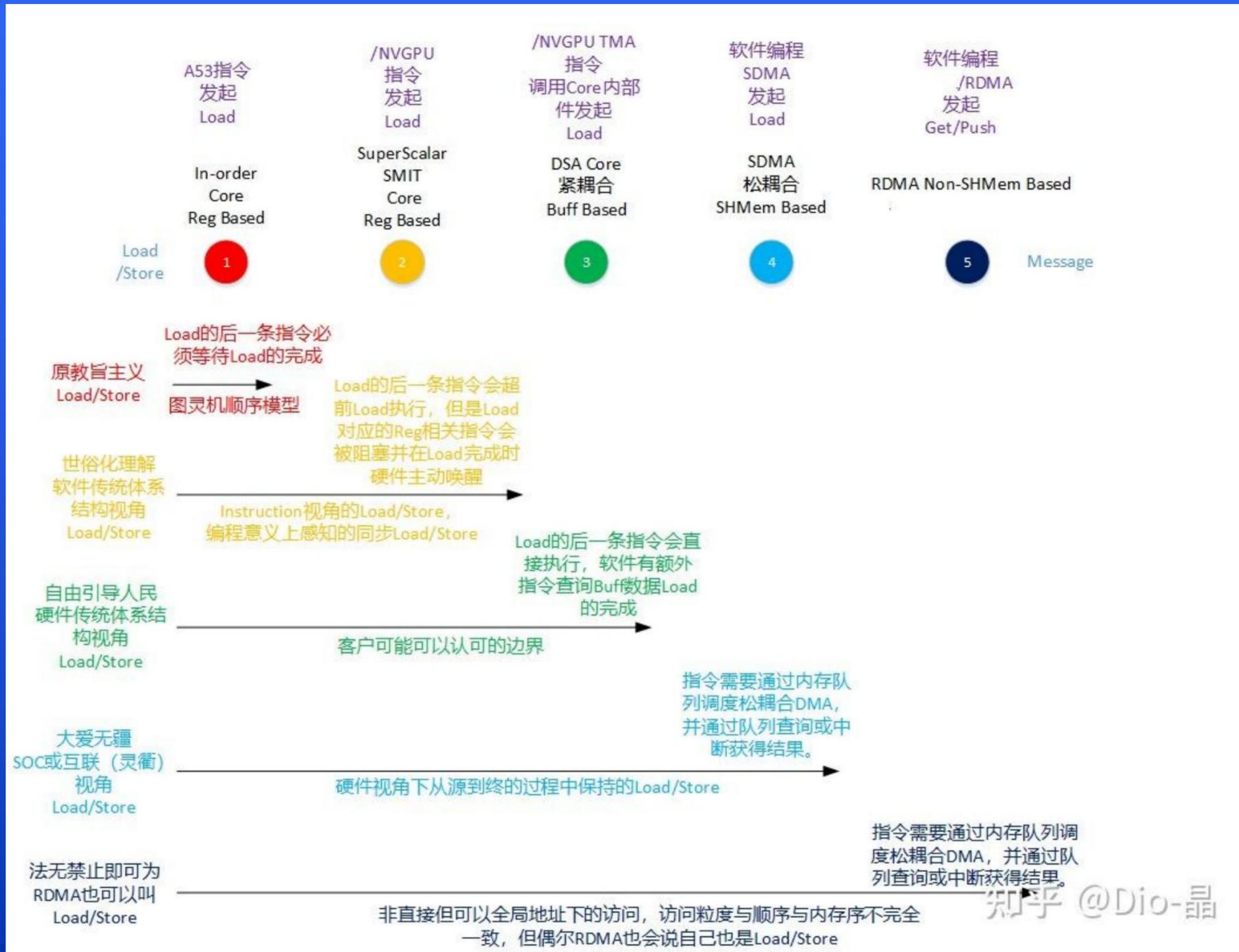


(b) UB-Mesh adopts the Unified Bus (UB) interconnects

CXL4.0 memory extension ?



内存语义 – Scale Up/Out 界限在哪里？



“

支持内存语义其实是指软件基于内存语义所表达出的丰富的软件行为，能够尽可能多且最低损耗地在Scale-Up系统上运行。

所以一个Scale-Up协议支持load/store这种原生内存访问行为是为了如上的“尽可能多”和“最低损耗”的达成，这是个度的概念，你可以无脑定个70%，或者换句话说，Load/Store更倾向于Native内存语义的表达。

为了保证更丰富和机制的内存语义在Scale-Up系统上无损的传递，一个好的Scale-Up协议应当不仅仅是Load/Store，还需要定制若干Macro行为来让某些更复杂内存语义在跨机时的高效表达，这是下一个哑谜，也是Ualink/Nvlink还没有做到的事情。

”

内存语义与编程范式讨论

内存语义 - 编程示例

```
1 #include <stdint.h>
2
3 // 功能: 将 src 地址的 len 字节数据拷贝到 dst 地址
4 void ldst_memcpy(uint8_t *dst,
5                 const uint8_t *src,
6                 uint32_t len) {
7     for (uint32_t i = 0; i < len; i++) {
8         // Load: CPU 从 src[i] 读取数据到寄存器
9         uint8_t data = src[i];
10        // Store: CPU 从寄存器将数据写入 dst[i]
11        dst[i] = data;
12    }
13 }
14
15 // 调用示例
16 int main() {
17     uint8_t src[1024] = "LD/ST Test Data";
18     uint8_t dst[1024] = {0};
19
20     ldst_memcpy(dst, src, sizeof(src)); // 执行拷贝
21     return 0;
22 }
```

```
1 #include <linux/dmaengine.h>
2 #include <linux/dma-mapping.h>
3
4 // 功能: 通过 DMA 控制器拷贝“内核内存 -> 外设缓冲区” (如网卡/SSD)
5 int dma_memcpy(struct device *dev, void *dst, const void *src, size_t len) {
6     struct dma_chan *dma_chan; // DMA 通道 (硬件分配)
7     struct dma_slave_config cfg; // DMA 配置参数
8     dma_addr_t src_dma, dst_dma; // DMA 可访问的物理地址 (非虚拟地址)
9     struct dma_async_tx_descriptor *desc; // DMA 传输描述符
10
11     // 1. 分配 DMA 通道 (根据设备类型, 如“内存-内存”或“内存-外设”)
12     dma_chan = dma_request_chan(dev, "dma0");
13     if (!dma_chan) return -ENODEV;
14
15     // 2. 映射虚拟地址到 DMA 物理地址 (让 DMA 硬件可见)
16     src_dma = dma_map_single(dev, (void *)src, len, DMA_TO_DEVICE); // 源: 内存-设备
17     dst_dma = dma_map_single(dev, dst, len, DMA_FROM_DEVICE); // 目标: 设备-内存
18     if (dma_mapping_error(dev, src_dma) || dma_mapping_error(dev, dst_dma))
19         return -EIO;
20
21     // 3. 配置 DMA 传输参数 (方向、地址、长度)
22     memset(&cfg, 0, sizeof(cfg));
23     cfg.direction = DMA_MEMCPY; // 传输方向
24     cfg.src_addr = src_dma; // 源 DMA 地址
25     cfg.dst_addr = dst_dma; // 目标 DMA 地址
26     cfg.src_addr_width = DMA_SLAVE_BUSWIDTH_4_BYTES; // 数据宽度 (4字节对齐)
27     dmaengine_slave_config(dma_chan, &cfg);
28
29     // 4. 创建并提交 DMA 传输任务
30     desc = dmaengine_prep_slave_single(
31         dma_chan, dst_dma, len, DMA_MEMCPY, DMA_PREP_INTERRUPT // 完成后触发中断
32     );
33     dmaengine_submit(desc);
34     dma_async_issue_pending(dma_chan); // 启动 DMA 传输
35
36     // 5. 等待 DMA 完成 (或通过中断回调处理)
37     dmaengine_wait_for_completion(dma_chan);
38
39     // 6. 释放资源
40     dma_unmap_single(dev, src_dma, len, DMA_TO_DEVICE);
41     dma_unmap_single(dev, dst_dma, len, DMA_FROM_DEVICE);
42     dma_release_chan(dma_chan);
43     return 0;
44 }
```

```
1 #include <infiniband/verbs.h>
2 #include <stdio.h>
3 #include <string.h>
4
5 --- write to remote side ---
6 #define REMOTE_IP "192.168.1.100" // 远程设备 IP
7 #define MEM_LEN 4096 // 传输数据长度 (4KB)
8
9 int main() {
10     struct ibv_context *ctx; // RDMA 设备上下文
11     struct ibv_pd *pd; // 保护域 (隔离内存区域)
12     struct ibv_mr *local_mr; // 本地内存区域 (MR)
13     char *local_buf; // 本地内存缓冲区
14     struct ibv_qp *qp; // 队列对 (用于 RDMA 通信)
15     struct ibv_send_wr wr = {0}; // 发送工作请求
16     struct ibv_sge sge = {0}; // 散射聚集元素 (描述内存区域)
17
18     // 1. 初始化 RDMA 设备 (获取上下文、保护域)
19     ctx = ibv_open_device_by_name(NULL, "mlx5_0"); // 假设使用 mlx5 网卡
20     pd = ibv_alloc_pd(ctx);
21
22     // 2. 分配本地内存并注册为 RDMA 内存区域 (MR)
23     local_buf = malloc(MEM_LEN);
24     strcpy(local_buf, "RDMA Test Data"); // 填充待传输数据
25     local_mr = ibv_reg_mr(
26         pd, local_buf, MEM_LEN,
27         IBV_ACCESS_LOCAL_WRITE | IBV_ACCESS_REMOTE_WRITE // 权限: 本地写、远程写
28     );
29
30     // 3. 与远端交换 MR 信息 (需通过 TCP 等协议, 此处简化)
31     // (实际需发送: local_mr->rkey (远程访问密钥)、local_mr->addr (物理地址))
32     struct {
33         uint64_t remote_addr; // 远程内存地址 (从远端获取)
34         uint32_t remote_rkey; // 远程内存访问密钥 (从远端获取)
35     } remote_info;
36     tcp_exchange_mr_info(REMOTE_IP, &remote_info, local_mr); // 自定义 TCP 交换
37
38     // 4. 配置 RDMA Write 工作请求 (写本地数据到远程内存)
39     sge.addr = (uint64_t)local_buf; // 本地内存地址
40     sge.length = MEM_LEN; // 传输长度
41     sge.lkey = local_mr->lkey; // 本地 MR 密钥
42
43     wr.opcode = IBV_WR_RDMA_WRITE; // 操作类型: RDMA 写
44     wr.sg_list = &sge; // 关联内存区域
45     wr.num_sge = 1;
46     wr.wr.rdma.remote_addr = remote_info.remote_addr; // 远程内存地址
47     wr.wr.rdma.rkey = remote_info.remote_rkey; // 远程 MR 密钥
48     wr.send_flags = IBV_SEND_SIGNALED; // 传输完成后发信号
49
50     // 5. 提交 RDMA Write 请求 (硬件直接执行)
51     ibv_post_send(qp, &wr, NULL);
52
53     // 6. 等待传输完成 (通过完成队列 CQ 检查)
54     wait_for_rdma_completion(qp); // 自定义等待函数
55
56     // 7. 释放资源
57     ibv_dereg_mr(local_mr);
58     free(local_buf);
59     ibv_dealloc_pd(pd);
60     ibv_close_device(ctx);
61     return 0;
62 }
63
64
65 ---- recv side ----
66
67 // 核心步骤: 分配远程内存 + 注册 MR + 通过 TCP 发送 MR 信息 (addr + rkey)
68 char *remote_buf = malloc(MEM_LEN);
69 struct ibv_mr *remote_mr = ibv_reg_mr(
70     pd, remote_buf, MEM_LEN,
71     IBV_ACCESS_REMOTE_WRITE // 允许远端读写此内存
72 );
73 tcp_send_mr_info(LOCAL_IP, remote_mr->addr, remote_mr->rkey); // 发送给本地端
```

SMP 兼容的共享内存 - 必要性?

```
11 // 定义共享内存的键值
12 #define SHM_KEY 1234
13
14 void *get_shm_ptr(void)
15 {
16     int shmid = shmget(SHM_KEY, sizeof(std::atomic<int>), IPC_CREAT | 0666);
17     if (shmid == -1) {
18         std::cerr << "Failed to create shared memory segment." << std::endl;
19         return 1;
20     }
21
22     void *shm_ptr = shmat(shmid, nullptr, 0);
23     if (shm_ptr == (void *)-1) {
24         std::cerr << "Failed to attach shared memory segment." << std::endl;
25         return 1;
26     }
27
28     return shm_ptr;
29 }
30
31 int main(int argc, char *argv[])
32 {
33     void *shm_ptr = get_shm_ptr();
34     #if 1
35         std::atomic<int> *g_atomic_var = new (shm_ptr) std::atomic<int>(0);
36     #else
37         std::atomic<int> *g_atomic_var = (std::atomic<int> *)shm_ptr;
38     #endif
39
40     std::cout << "g_atomic_var value = "
41               << g_atomic_var->load() << std::endl;
42     g_atomic_var->fetch_add(5);
43     std::cout << "g_atomic_var value = "
44               << g_atomic_var->load() << std::endl;
45
46     g_atomic_var->fetch_sub(2);
47     std::cout << "g_atomic_var value = "
48               << g_atomic_var->load() << std::endl;
49     auto old_value = g_atomic_var->exchange(10);
50     std::cout << "old_valud: " << old_value
51               << ", g_atomic_var value = " << g_atomic_var->load()
52               << std::endl;
53
54     auto r2 = g_atomic_var->compare_exchange_strong(old_value, 11);
55     std::cout << "r2: " << r2
56               << ", g_atomic_var value = " << g_atomic_var->load()
57               << std::endl;
58
59     shmdt(shm_ptr);
60     shmctl(shmid, IPC_RMID, nullptr);
61     return 0;
62 }
```

可换为: get_global_shm_ptr

Transaction based

SMP style Lock based

```
6 static auto v0 = 0, v1 = 0;
7 std::mutex g_lock;
8
9 void swap(int& a, int& b) {
10     int temp = a;
11     a = b;
12     b = temp;
13 }
14
15 // 事务函数, 用于交换两个整数的值
16 void swap_a(int& a, int& b) __transaction_atomic {
17     int temp = a;
18     a = b;
19     b = temp;
20 }
21
22 static void thread_f(void)
23 {
24     for (auto i = 0; i < 1024 * 1024; i++) {
25         #if defined(__TA)
26             __transaction_atomic
27         #else
28             std::lock_guard<std::mutex> lock(g_lock);
29         #endif
30         {
31             v0++;
32             swap(v0, v1);
33             v1++;
34             swap(v1, v0);
35         }
36     }
37 }
38
39 int main() {
40     int x = 10;
41     int y = 20;
42
43     std::cout << "Before swap: x = " << x << ", y = " << y << std::endl;
44
45     // 调用事务函数进行交换
46     swap_a(x, y);
47
48     std::cout << "After swap: x = " << x << ", y = " << y << std::endl;
49
50     std::cout << "Before tp: v0 = " << v0 << ", v1 = " << v1 << std::endl;
51
52     std::thread tp[16];
53     for (auto &t : tp) {
54         t = std::thread(thread_f);
55     }
56
57     for (auto &t : tp) {
58         t.join();
59     }
60
61     std::cout << "After tp: v0 = " << v0 << ", v1 = " << v1 << std::endl;
62     return 0;
63 }
```

亦可为: 其他 Bus Native Semantics - TBD

Scale-UP 计算系统中的内存模型 – 如何定义?

比如NVIDIA GPU 的Memory Consistency Model 定义了三种保序模式:

(1) **Acquire**: 标记为Acquire类型的Load Store, 可以看作一个前fence操作, 后面的relaxed 类型或者Acquire类型的指令不能被乱序到前面执行。

Acquire

```
cuda::atomic<int> a;

//Prior load
int before = array[0];

//Atomic load
val = a.load(std::memory_order_acquire);

// Later load
int after = array[0];
```

- Loads and stores **cannot** be moved **before** the acquire
- **Same-address ordering** is preserved within a single thread

(2) **Release**: 标记为Release类型的Load Store, 可以看作一个后fence操作, 写在前面的relaxed 类型或者Release类型的指令不能被乱序到后面执行。

Release

```
cuda::atomic<int> a;

//Prior load
int before = array[0];

//Atomic load
val = a.load(std::memory_order_release);

// Later load
int after = array[0];
```

- Loads and stores **cannot** be moved **after** the release
- **Same-address ordering** is preserved within a single thread

(3) **Relaxed**: 标记为Relaxed类型的Load/Store 之间可以乱序执行

Relaxed

```
cuda::atomic<int> a;

//Prior load
int before = array[0];

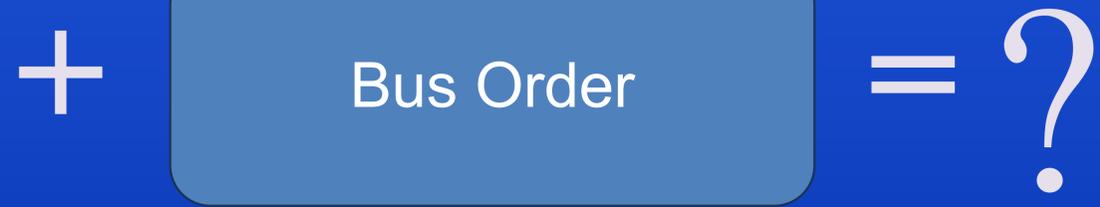
//Atomic load
val = a.load(std::memory_order_relaxed);

// Later load
int after = array[0];
```

- Loads and stores **can** be moved **before** and **after** relaxed
- **Same-address ordering** is preserved within a single thread

Memory ordering in some architectures^{[8][9]}

Type	Alpha	ARMv7	MIPS	RISC-V		PA-RISC	POWER	SPARC			x86 [a]	AMD64	IA-64	z/Architecture
				WMO	TSO			RMO	PSO	TSO				
Loads can be reordered after loads	Y	Y		Y		Y	Y	Y					Y	
Loads can be reordered after stores	Y	Y		Y		Y	Y	Y					Y	
Stores can be reordered after stores	Y	Y		Y		Y	Y	Y	Y				Y	
Stores can be reordered after loads	Y	Y	depend on implementation	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y
Atomic can be reordered with loads	Y	Y		Y			Y	Y					Y	
Atomic can be reordered with stores	Y	Y		Y			Y	Y	Y				Y	
Dependent loads can be reordered	Y													
Incoherent instruction cache pipeline	Y	Y		Y			Y	Y	Y	Y	Y		Y	



> <https://zhuanlan.zhihu.com/p/1922796547814458461>

讨论 & 展望

THANKS