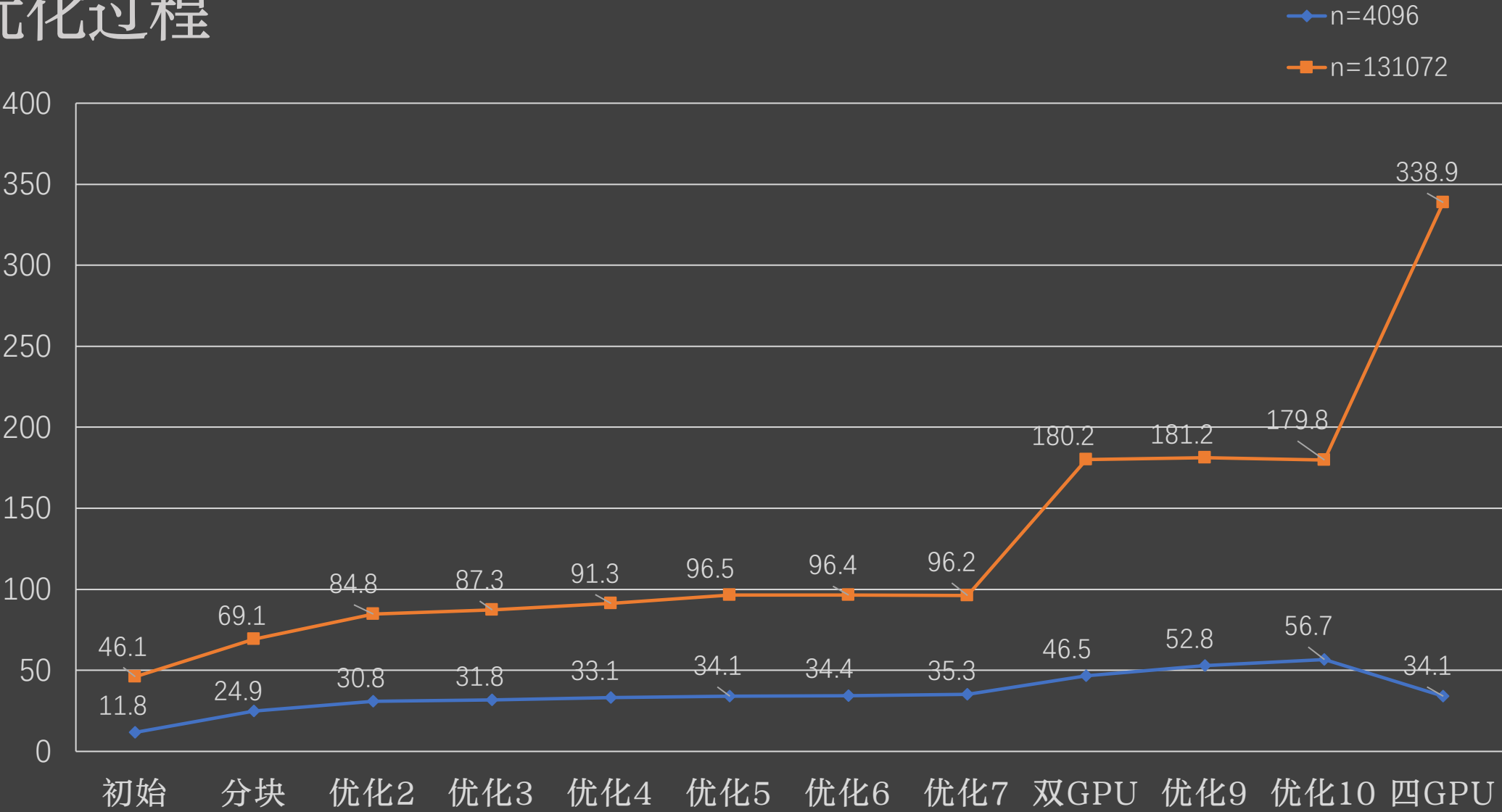


CUDA N-body优化

By 刘畅

优化过程



1

12→25 ($n = 2^{12}$)

46→69 ($n = 2^{17}$)

利用Shared Memory分块计算

» 分块大小等于block大小

```
__global__ void
bodyUpdateVelocity(float3 *pos, float3 *vel, float dt, int n) {
    int bodyId = blockIdx.x * blockDim.x + threadIdx.x;
    if (bodyId >= n) return;
    extern __shared__ float3 posShared[];

    float3 acc = { 0.0f, 0.0f, 0.0f };
    float3 p = pos[bodyId];

    for (int i = 0; i < gridDim.x; i++) {
        int id = i * blockDim.x + threadIdx.x;
        __syncthreads();
        posShared[threadIdx.x] = pos[id];
        __syncthreads();
        bodyGatherAccelTiled(p, acc);
    }
    vel[bodyId] += acc * dt;
}
```

2

25→31

69→85

循环展开

- » 大概还是由于降低了循环中分支的频率
- » #pragma unroll自动预处理

```
__forceinline__ __device__  
void bodyGatherAccelTiled(float3 pos, float3 &acc) {  
    extern __shared__ float3 posShared[];  
    #pragma unroll 32  
    for (int i = 0; i < blockDim.x; i += 2) {  
        bodyGetAccelPairwise(pos, posShared[i + 1], acc);  
        bodyGetAccelPairwise(pos, posShared[i], acc);  
    }  
}
```

⋈再尝试手动展开一层，便于调整shared mem访问顺序（微微作用）

3

34.1→34.4

96.5→96.4

CUDA PTX内联汇编优化

» 乘加指令： $D = A \times B + C$

比较明显的模式编译器能够自动识别，不过有例外：

`distSqr = dx * dx + dy * dy + dz * dz + SOFTENING`

```
__forceinline__ __device__
float dot(const float3 &a, const float3 &b) {
    const float SOFTENING = 1e-9f;
    float res = SOFTENING;
    asm("fma.rn.f32 %0, %1, %2, %3;" :
        "=f"(res) : "f"(a.x), "f"(b.x), "f"(res));
    asm("fma.rn.f32 %0, %1, %2, %3;" :
        "=f"(res) : "f"(a.y), "f"(b.y), "f"(res));
    asm("fma.rn.f32 %0, %1, %2, %3;" :
        "=f"(res) : "f"(a.z), "f"(b.z), "f"(res));
    return res;
}
```

» 效果极其有限

双GPU，双倍快乐

» 两张Tesla K80可用，然而

GPU Fan	Name Temp Perf	Persistence-M Pwr:Usage/Cap	Bus-Id	Disp.A Memory-Usage	Volatile GPU-Util	Uncorr. Compute M. ECC
0 N/A	Tesla K80 38C P8	Off 27W / 149W	00000000:09:00.0	Off 11MiB / 12206MiB	0%	Off Default
1 N/A	Tesla K80 32C P8	Off 28W / 149W	00000000:0A:00.0	Off 11MiB / 12206MiB	0%	Off Default
2 N/A	Tesla K80 59C P0	Off 57W / 149W	00000000:86:00.0	Off 582MiB / 12206MiB	30%	Off Default
3 N/A	Tesla K80 43C P0	Off 70W / 149W	00000000:87:00.0	Off 321MiB / 12206MiB	0%	Off Default

» 原来一张卡上有两个

NVIDIA Tesla K80

GK210 x2 GRAPHICS PROCESSOR	2496 x2 CORES	208 x2 TMUS	48 x2 ROPS	12 GB x2 MEMORY SIZE	GDDR5 MEMORY TYPE	384 bit x2 BUS WIDTH
--------------------------------	------------------	----------------	---------------	-------------------------	----------------------	-------------------------

» 先试试二路

4

35→46

96→180

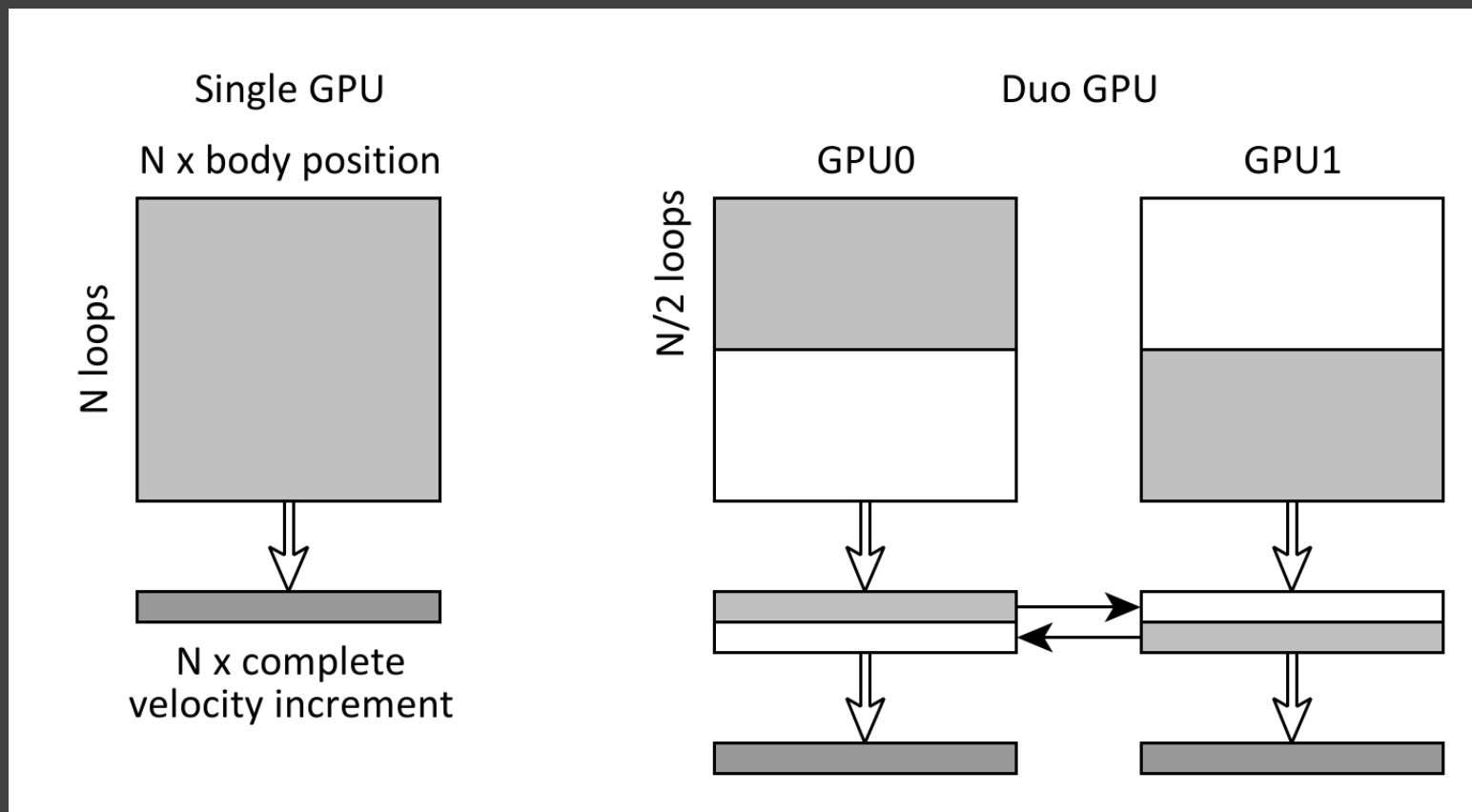
双GPU，双倍快乐

- » 如何划分 n^2 对加速度的计算？
- » GPU数据怎样共享/传输/同步？

4

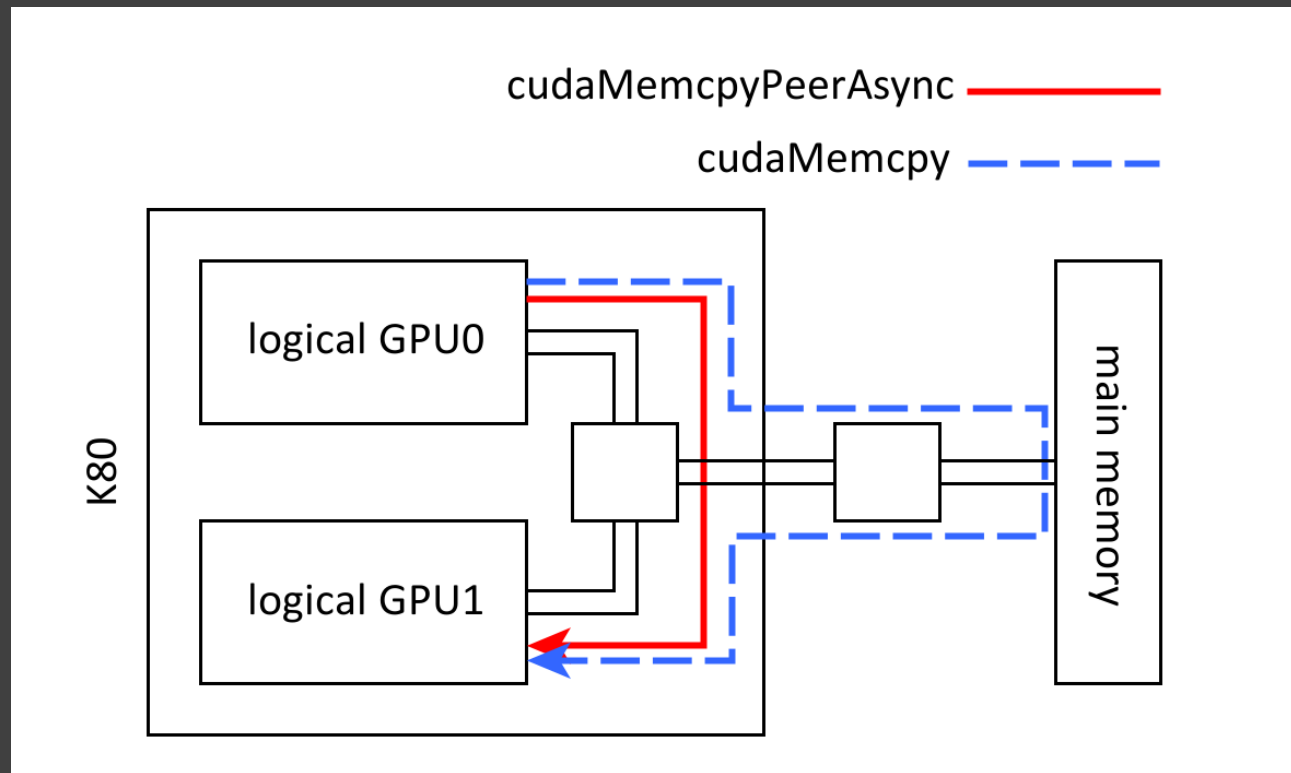
35→46

96→180



优化数据传输

» 直接走内部总线



» 需要peer access支持

5

46→53

180→181

继续利用peer access

6

53→57

181→180

» `cudaDeviceEnablePeerAccess`

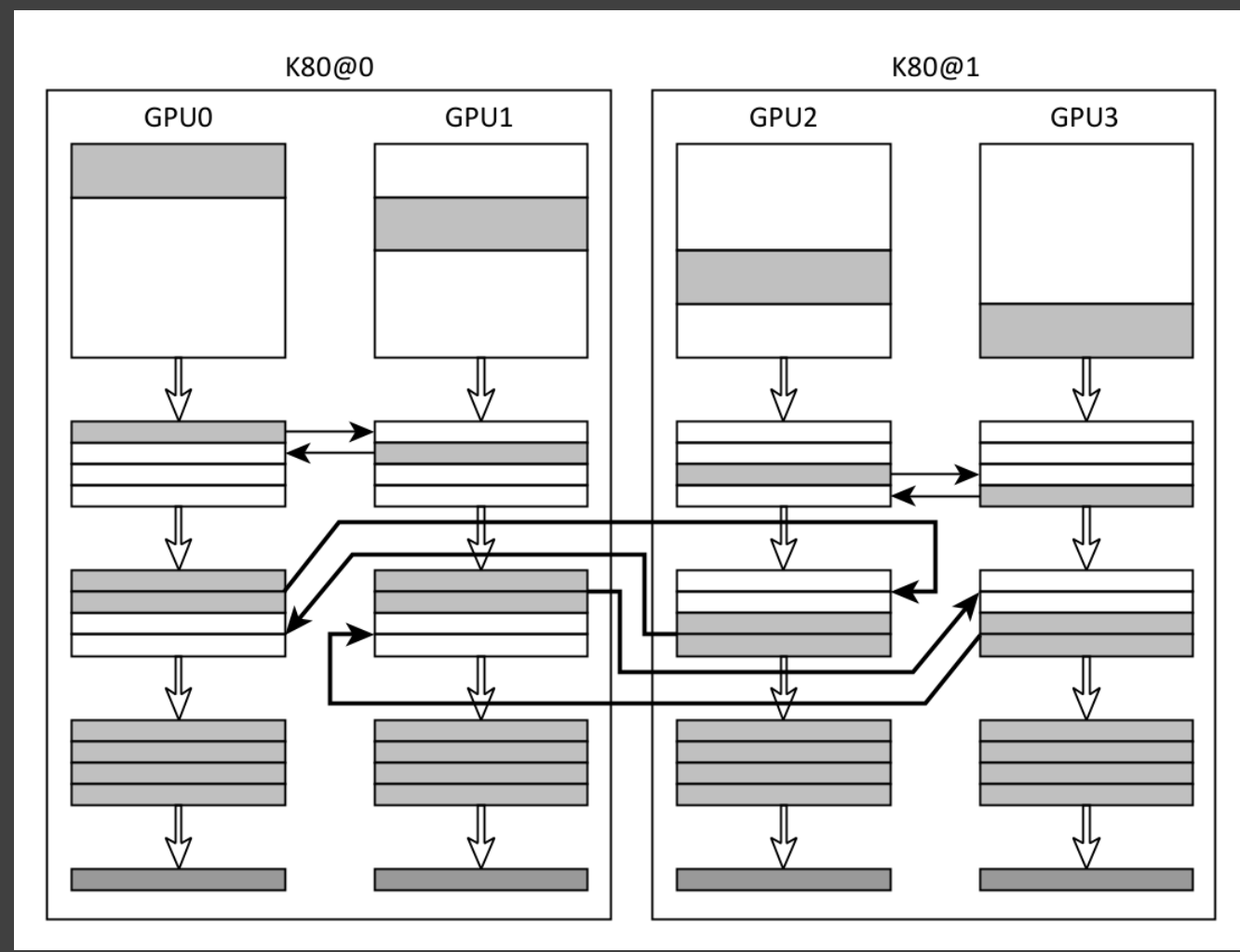
- 这样两个GPU可以直接访问对方的内存地址空间
- 直接累加对方的速度增量到自己算出来的部分
- 不稳定

双双GPU，四倍快乐

7

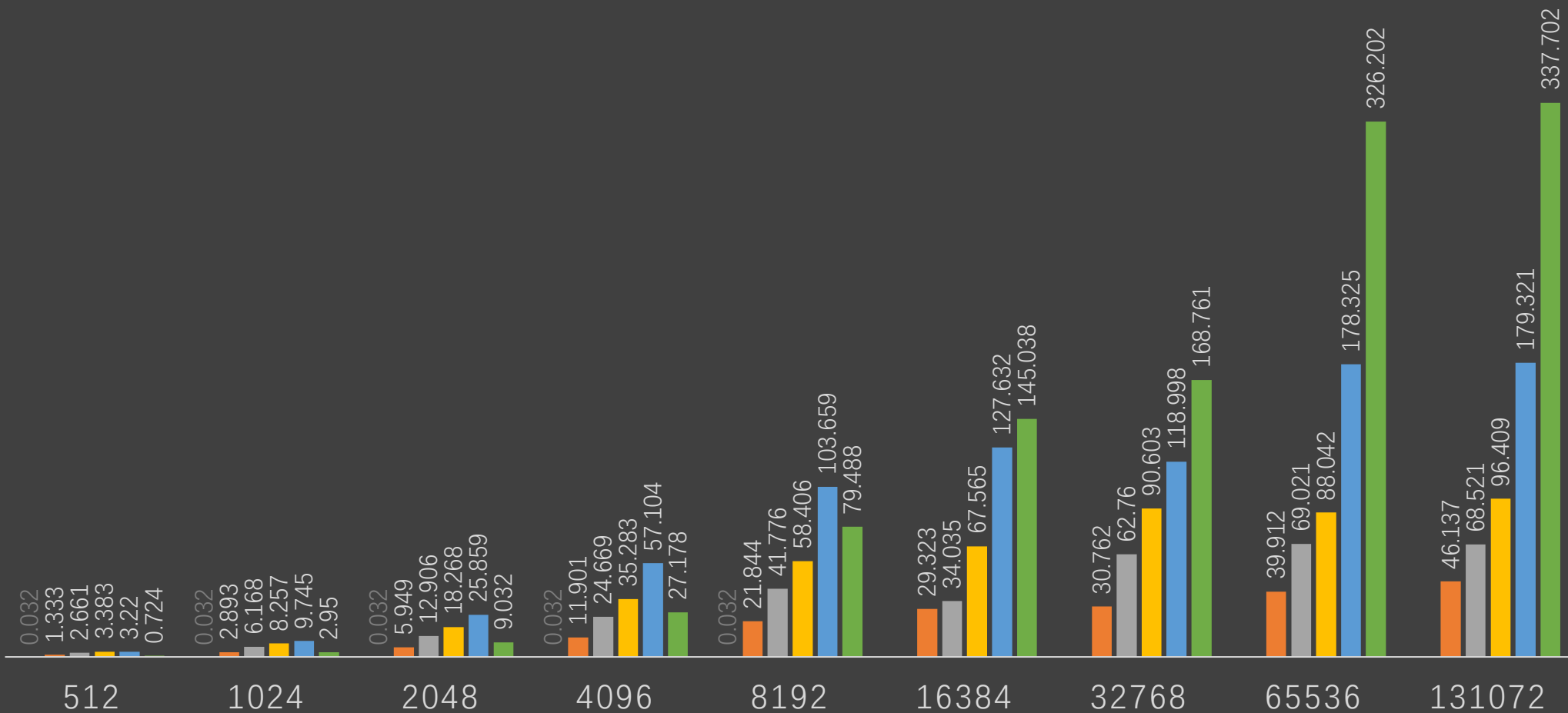
57→34

180→339



优化对比

■ 串行 ■ 并行 ■ 分块 ■ 分块+优化 ■ 双GPU ■ 4GPU



其他方法

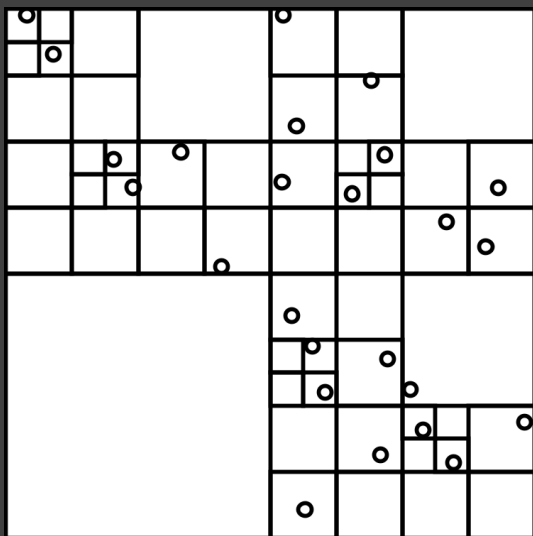
2. 基于空间划分和远距离近似的快速算法（硬核）

» 基于以下近似假设

当多个body聚集在很小一团时，其对很远之外的一个body的作用可近似为单个整体

» 因此，可以将body按空间划分

对于离得远的那些body，就不必逐一计算



1. 多thread对一body

» 小规模数据较快

» 大规模数据吞吐率提升有限

可能是由于大量原子操作导致的串行化

» 代表性算法

· Barnes-Hut $O(n \log n)$

· Fast Multipole Method $O(n)$

» 每次计算前都要并行构建八叉树 ←KD树和BVH应该也可以
或者动态调整结构

» 结点需要保存质量、质心、包围盒（球）信息

» 当结点内的多个body的包围盒相对被计算body的立体角小于某个阈值，可近似为单个质点