

--shfl\_down\_sync(mask, val, offset)

① participating thread calling this function.

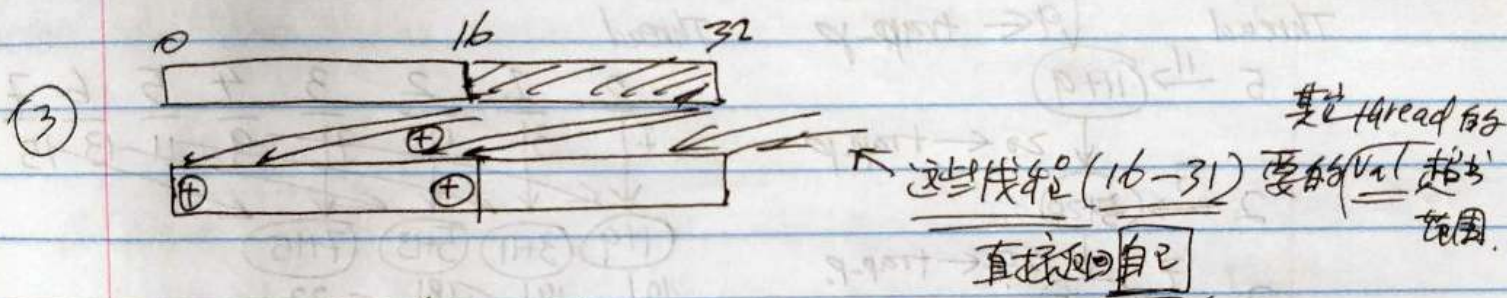
if mask 对应那个线程  $\Rightarrow$  这个线程从别的线程处获取 val, 是吗?  
 为 0 就得直接返回自己?

但是 doc 说  
 undefined  
 behavior

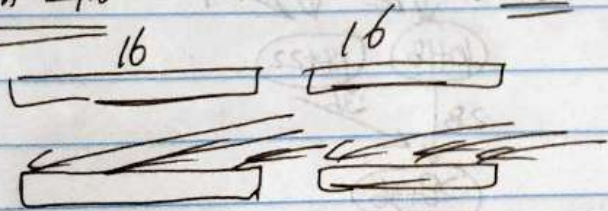
实验来看, mask = 0. (全 0)  
 调用线程仍然工作 返回对应线程的值.

② mask 为 0xffffffff (全 1), 但是有 31 个线程不调用这个函数 (inactive)  
 那些线程去拿它的值后返回 是不确定的

因为那些线程并没有提供值去 shuffle



④ width = 16 16 个 thread 为一组, 不可跨 2 个 warp.



超出范围, down\_sync (返回自己).

shfl\_sync 返回 wrap up 结果  $\frac{srcLane \% width}{width}$

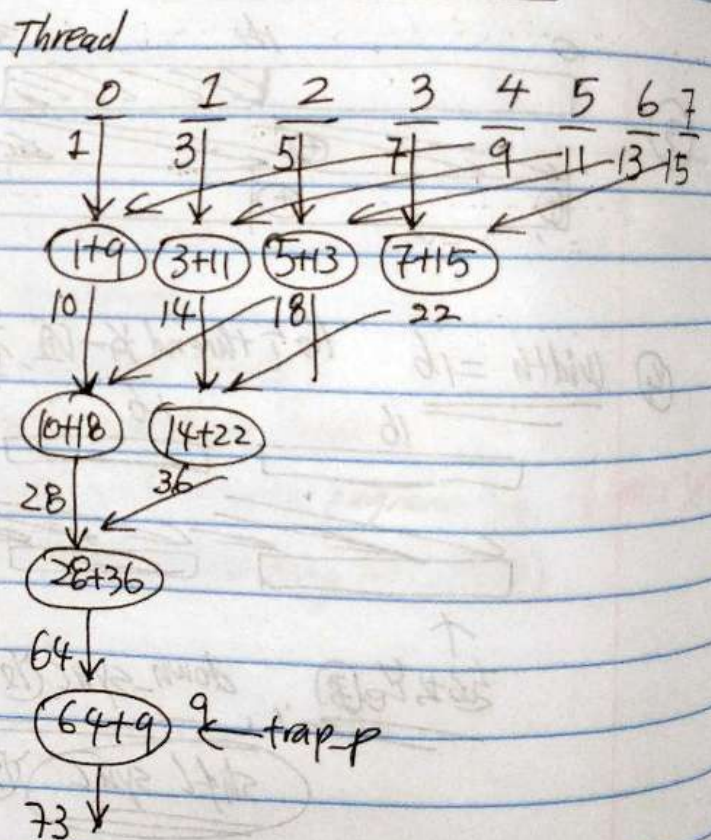
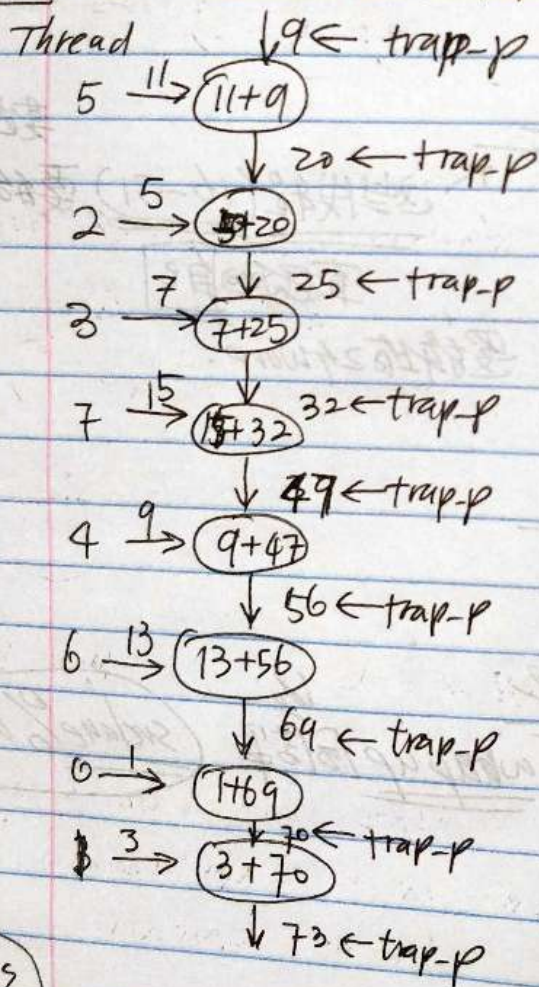


# SIMT: Single instruction Multiple Thread (not SIMD)

because threads on an SM that are executing the same instruction may not execute simultaneously; to hide memory access latency, some threads may block while memory is accessed and other threads, that have already accessed the data, may proceed with execution.

SIMD: composed of a single control unit and multiple datapaths. The control unit fetches an instruction from memory and broadcast it to the datapaths. Each data path either executes the instruction on its data or is idle.

Basic Sum  $\text{old value} \leftarrow \text{AtomicAdd}(\text{trap-p}, \text{val});$  Tree-structured Sum



$t_{\text{threads}} + t_{\text{value}} \rightarrow t_{\text{sequential additions}}$

$$\lceil \log_2 t \rceil + 1$$



warpSize

$$\text{lane} = \text{threadIdx.x} \% \text{warpSize};$$

the warp shuffle function allows threads in a warp to read from registers used by another thread in the same warp.

--device-- float --shfl\_down\_sync(  
    unsigned mask,  
    float var,  
    unsigned diff, // l+diff (diff > 0).  
    int width = warpSize).

A bit representing the thread's lane, must be set for each participating thread to ensure that all of the threads in the call have converged - i.e. arrived at the call before any thread begins executing the call to --shfl\_down\_sync.

thread l      value  
            ← thread l+diff

① what happens if thread l calls --shfl\_down\_sync, but l+diff doesn't??

↳ the value returned by the call on thread l is undefined (l+diff thread is inactive due to if branch) ↓ 返回 0 或 NaN

② what happens if thread l calls --shfl\_down\_sync but l+diff > warpSize

↳ the call returns the value in var already stored on thread l.  
(original value of var) ② width 控制

③ what happens if thread l calls --shfl\_down\_sync but l+diff < warpSize, and l+diff > largest lane in the warp. (due to ~~warp~~ thread block size is not multiplier of warp size). the last warp has no 32 threads (< 32).

similar to ①, some threads are inactive.

so value returned by the call is also undefined 返回 0 或 NaN

mask 指定 participating threads 返回 ... (地址描述)

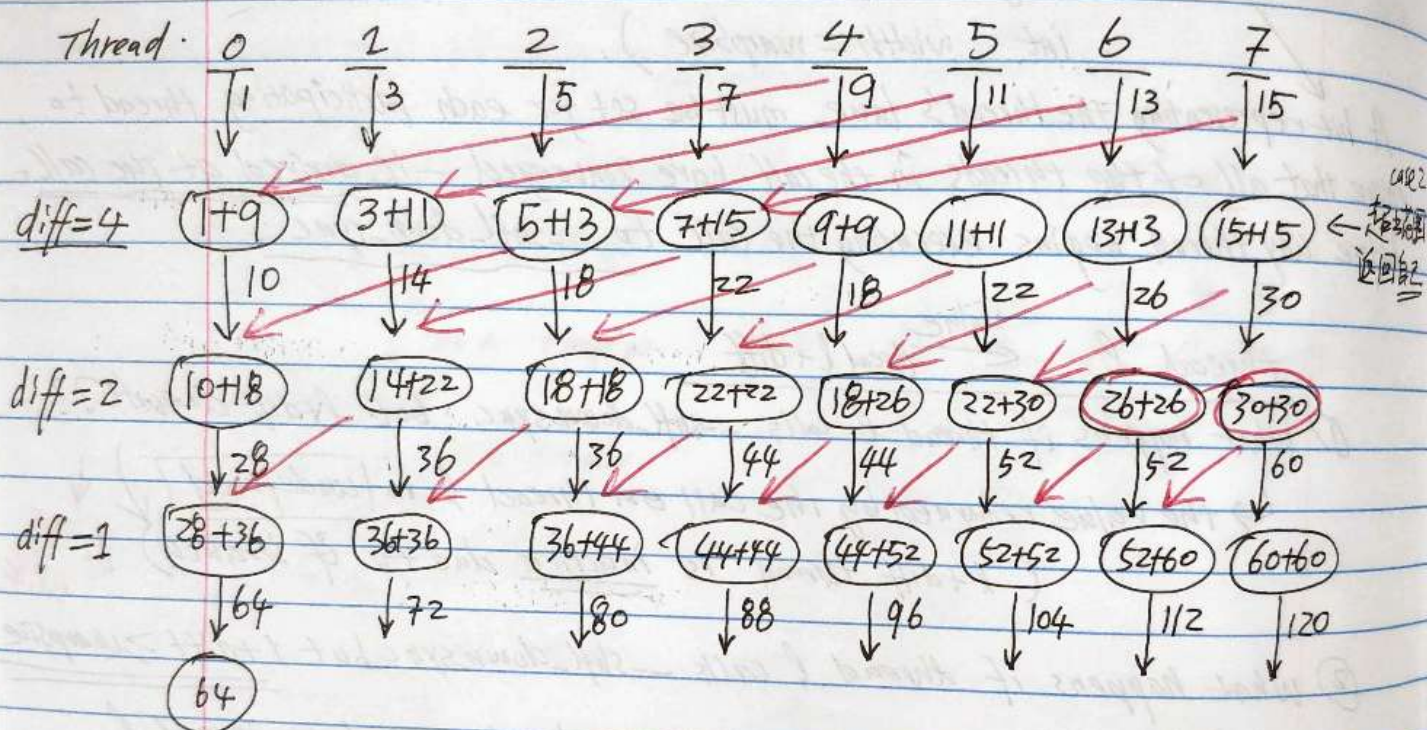
mask 指定 participating threads 返回 0 (inactive threads, ② mask 指定为 0).



```

--device_ float warpSum(float val) {
    unsigned mask = 0xffffffff;
    for (int diff = warpSize/2; diff > 0; diff >>= 1) {
        val += shfl_down_sync(mask, val, diff);
    }
    return val;
}

```



```

int my_i = blockDim.x * blockIdx.x + threadIdx.x;

```

```

if (float my_trap = 0.0f;

```

```

    if (0 < my_i && my_i < n) {

```

```

        float my_x = a + my_i * h;

```

```

        my_trap = f(my_x);
    }

```

```

float result = warpSum(my_trap);

```

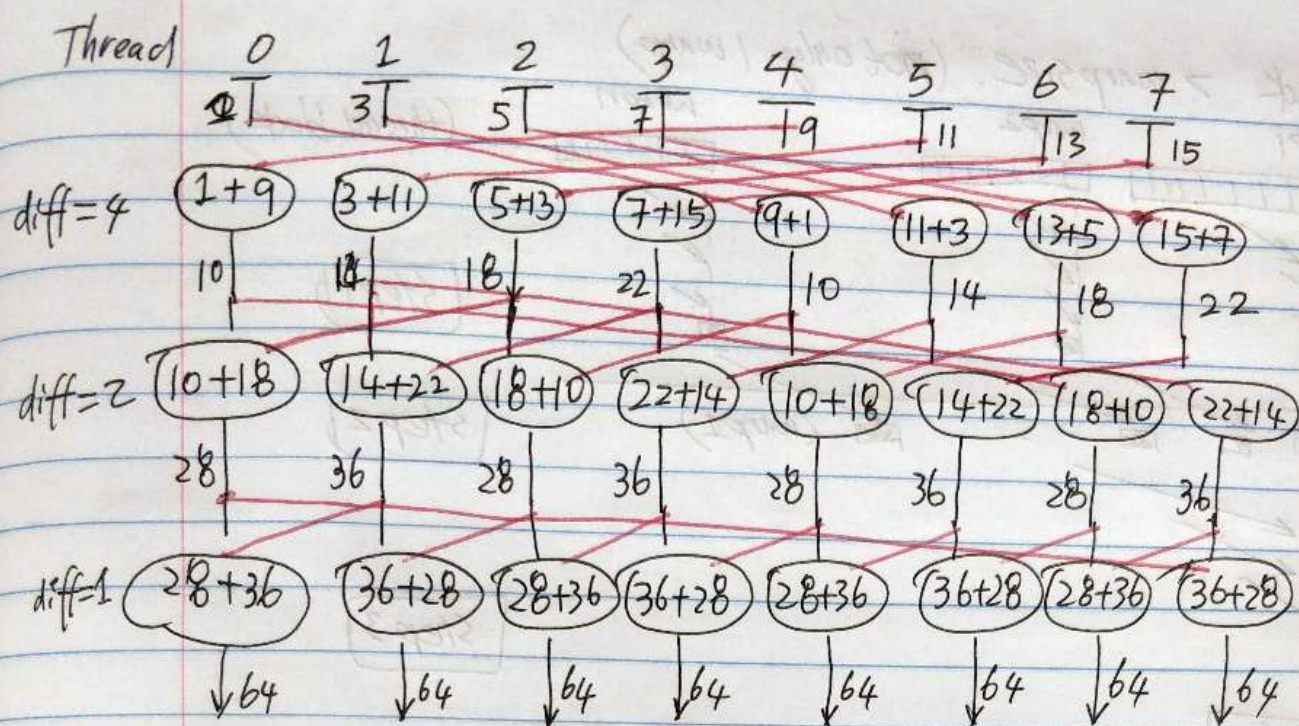
```

if (threadIdx.x == 0) atomicAdd(&trap_p, result);

```

(one thread block only a warp)





← every thread will have final sum.

```

__device__ float shared_mem_sum(float shared_vals[7]) {
    int my_lane = threadIdx.x % warp_size;
    for (int diff = warp_size / 2; diff > 0; diff >>= 1) {
        int source = (my_lane + diff) % warp_size;
        shared_vals[my_lane] += shared_vals[source];
    }
    return shared_vals[my_lane];
}

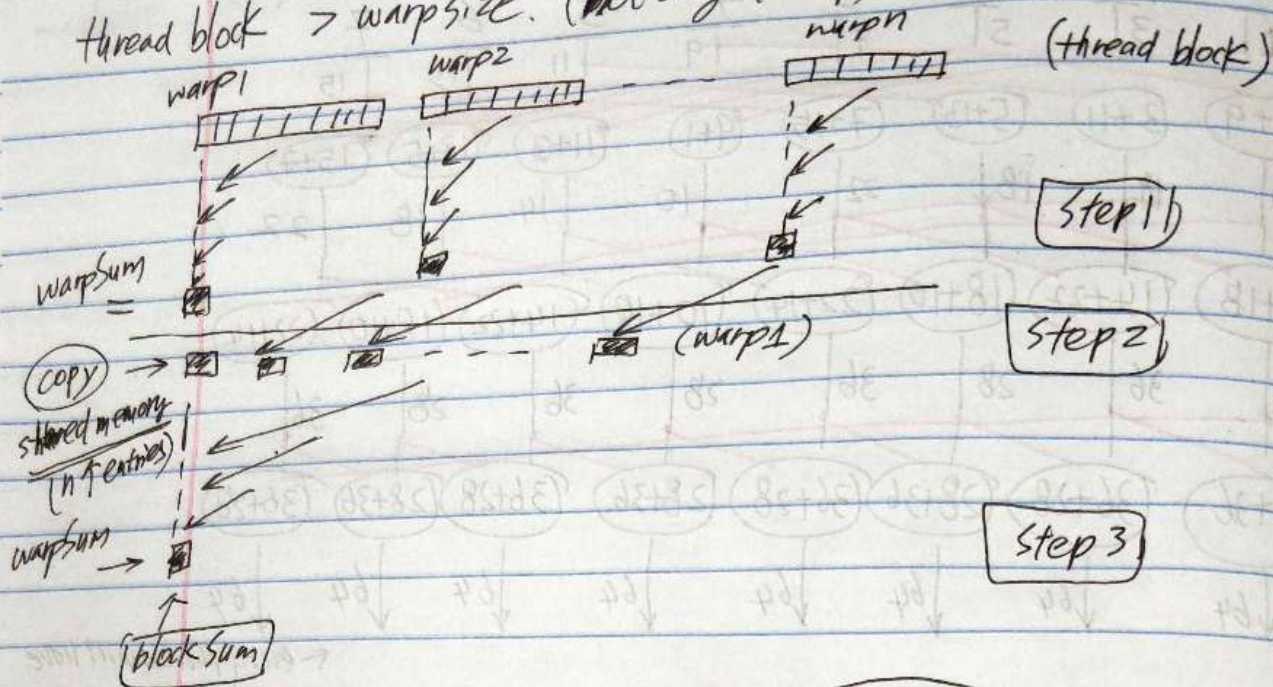
```

// dissemination sum using shared memory

← either global memory or shared memory



thread block  $\rightarrow$  warpsize. (not only 1 warp)



pytorch code

```
T BlockReduceSum(T val, T* shared) {
```

```
    const int laneid = threadIdx.x % warpsize;
```

```
    const int warpid = threadIdx.x / warpsize;
```

```
    val = WarpReduceSum(val);
```

```
    __syncthreads(); // 线程同步!!!
```

```
    if (laneid == 0)
```

```
        shared[warpid] = val;
```

```
    __syncthreads();
```

```
    val = (threadIdx.x < blockDim.x / warpsize) ? shared[laneid] : T(0);
```

```
    if (warpid == 0)
```

```
        val = WarpReduceSum(val);
```

```
    return val;
```

```
int sum = 0;
```

```
for (int i = blockIdx.x * blockDim.x + threadIdx.x;
```

```
    i < N;
```

```
    i += blockDim.x * gridDim.x) {
```

```
    sum += in[i];
```

```
}
```

```
sum = blockReduceSum(sum);
```

```
if (threadIdx.x == 0)
```

```
    atomicAdd(out, sum);
```

```
int sum = 0;
```

```
for (int i = blockIdx.x * blockDim.x + threadIdx.x;
```

```
    i < N;
```

```
    i += blockDim.x * gridDim.x) {
```

```
    sum += in[i];
```

```
}
```

```
sum = warpReduceSum(sum);
```

```
if (threadIdx.x & (WARP_SIZE - 1) == 0)
```

```
    atomicAdd(out, sum);
```

warps

32 32 32 ... 32

32 threads per warp. 11 warps

how many warps in thread block

↑ ↑

(warp)

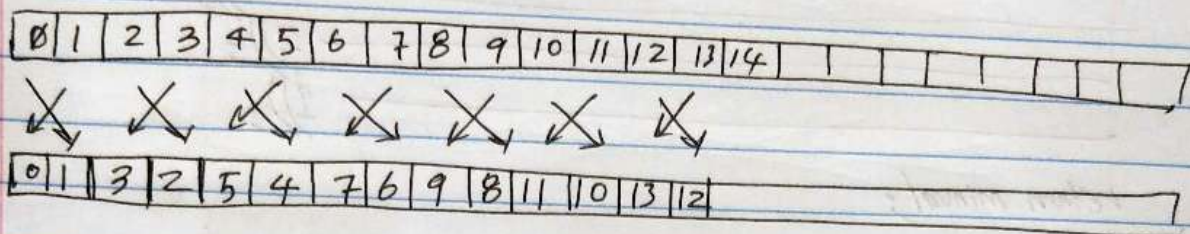
Also good.

skip block-level processing

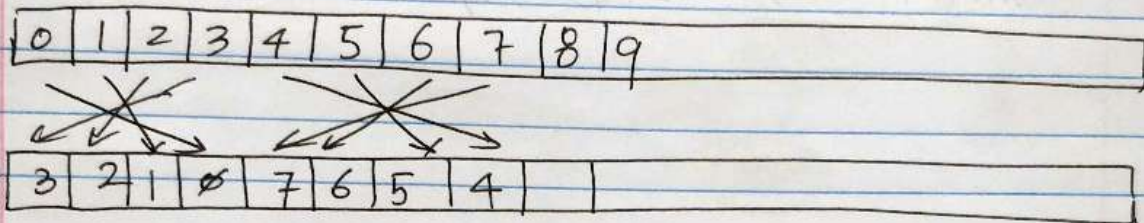


--shfl\_xor\_sync (unsigned mask, T var, int laneMask, int width=WarpSize)

laneMask=1

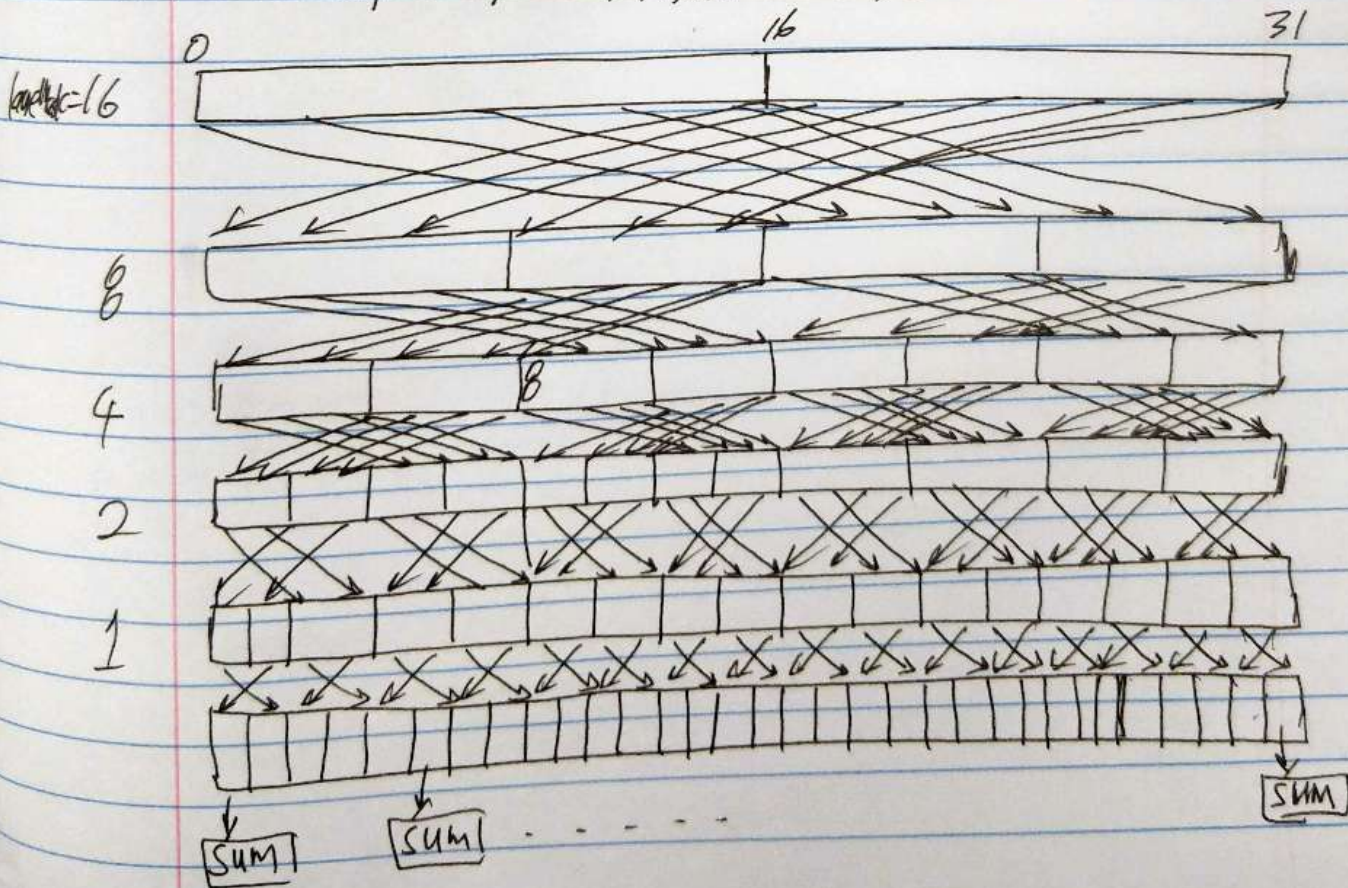


laneMask=3



`sum += --shfl_xor_sync(0xffffffff, sum, 16);`  
`sum += --shfl_xor_sync(0xffffffff, sum, 8);`  
`sum += --shfl_xor_sync(0xffffffff, sum, 4);`  
`sum += --shfl_xor_sync(0xffffffff, sum, 2);`  
`sum += --shfl_xor_sync(0xffffffff, sum, 1);`

也利用 ~~快速~~ <sup>min</sup> product;  
所有线程最后都得到 value  
meaningful





```

--device-- --forceline-- int warpReduceMin(int minVal)
minVal = min(minVal, --shfl_xor_sync(0xffffffff, minVal, 16))
minVal = min(minVal, --shfl_xor_sync(_____, _____, 8));
_____)4);
_____)2);
_____)1);
return minVal;
}

```

min  $\Rightarrow$  minVal  $\ast =$  --shfl\_xor\_sync \_\_\_\_\_ 16  $\rightarrow$  8  $\rightarrow$  4  $\rightarrow$  2  $\rightarrow$  1