



F 0  
 T F  
 T  
 T F 0  
 T F  
 T

~~~~~ throughput  
 O  
 ~~~~~ → SM, wrap, cores

Faster accessing read-only data  
~~wasting memory space~~  
 limited size restricts large size vectors

~~it takes 6 sec for the 10% workload cannot offload~~

~~→ 34 sec in CPU, 4 sec in GPU to CPU~~

~~⇒  $\frac{34}{4} = 8.5$  times speed up~~

$$40 \times 0.9 = 36 \text{ Sec} \rightarrow 4 \text{ sec}$$

$$40 \times 0.1 = 4 \text{ sec (CPU)} \quad \downarrow$$

9.1倍

$$10 - 4 - 4 = 2 \quad \#$$

~~~~~

$$\boxed{3 \times 3 \times 3}$$

$$5 \times (32-2)(32-2) \times [2 \times 2]$$

$$(32-2)(32-2)$$

✓ a ,

→ improving cache

✓ c<sup>0</sup> //

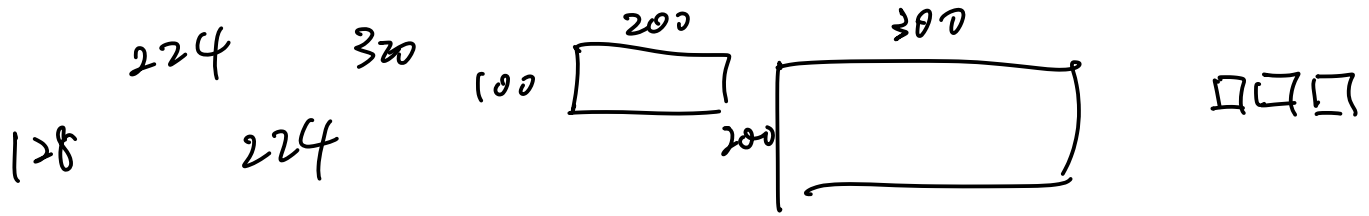
✓ b

→ threads within same wrap will access consecutive global memory

✓ b o

b

✓



Ask does halo cells require FLOP?

○

$$(200 \times 2) \times 100 \times 300$$

or  $(224 \times 2) \times (128 \times 320)$ ?

$$\frac{300}{32} \approx 10$$

$$10 \times 200 \times 100 + 4 \times 200 \times 300$$

$$\frac{8 \times 10 \times 12}{8 \times 8 \times 8}$$

$$\frac{8 \times 10 \times 12}{8 \times 8 \times 8} = 1.875$$

$$\frac{8 \times 8 \times 8 \times 1 \times 3 \times 5}{8 \times 10 \times 12} = 8$$

○ threadIdx.x 0 ~ 127, 128 ~ 255 are in same wrap.

~~1024~~ 128

—

\_\_\_\_\_

~~~~~

0 global ✓ shared ✓ register ✓ register ✓

$$\text{int index} = \text{blockId.x} \cdot \text{BlockElement} \\ + i \cdot \text{ElemThread} + \text{threadId.x}$$

⇒ use coalescing to load consecutive global addresses


Control divergence occurs when threads within the same warp take different execution paths due to conditional branches. In such cases, the GPU serializes the execution of the divergent branches, reducing parallel efficiency.

0 0

$$\cancel{32} \cdot 0 \\ \cancel{0} \cdot \frac{128}{32} = 4$$

0

The error arises because multiple threads update the shared variable max concurrently without proper synchronization or atomic operations, causing race conditions.

  
⇒ tile = 32  
y  
-  
e

~~3~~ 32



```
-- shared -- tile[TILE_WIDTH][tile_width]
int count = 0
```

```
-
tile[tz][ty] = in[index_z * y_size + index_y]
- synthreads()
```

```
z < index_z + blur_r
y < index_y + blur_y
z < 0 or z > row or y < 0 or y > col
```

if in current block

~~$0 \leq tz \leq \text{tile\_width}$~~

~~$0 \leq ty \leq \text{tile\_width}$~~

~~$z - bz \times \text{block\_dim}_z$~~

~~$ty - blur_y$~~

$\text{sum} += \text{in}[z * y\_size + y]$

if (count > 0 && index\_z < z\_size && index\_y < y\_size)

$\text{out}[\text{index\_z} * y\_size + \text{index\_y}] = \text{sum} / \text{count}$

```

float *in_d
float *out_d
cudaMalloc((void**) &in_d, row * col * sizeof(float))
cudaMemcpy(in_d, in, row * col * sizeof(float), Host to Device)
dim3 dimGrid = (1, ceil(col/tilewidth), ceil(row/tilewidth))
dim3 dimBlock = (1, tilewidth, tilewidth)
blurKernel <<< dimGrid, dimBlock >>> (in_d, out_d,
cudaMemcpy(out, out_d, size, Device to Host) row, col, h(r, r)
cudaFree

```

Make sure in elements are completely loaded to the tile.

②

