

$[B, A] \rightarrow [A]$

## 1. Warp Read Warp Reduce

### \* Basic Method

- ① Each warp reads 32 consecutive values.  
A block reads 32 lines. ? ( $32 \times 32$  or  $16 \times 64$ )
- ② Save values in shared memory.  
Each warp reads a column.
- ③ Use warp reduce.

### \* Requirement for Efficiency

$$B \geq 32$$

### \* Different cases

- ①  $\lfloor A/32 \rfloor \geq N_{sm}$   
Each warp reads multiple lines and reduce first.

- ②  $\lfloor A/32 \rfloor < N_{sm}$   
Each warp reads 1 line.  
Different blocks reduce different part of the same column.  
Use atomic add for final result.

- ③ combination of ① and ②

### \* Corner Case / Margin Part ( $a = A < 32$ or $a = A \bmod 32$ )

$$\text{let } n = \lfloor 32/a \rfloor$$

- ① 32 warps, each warp reads  $n$  lines (margin part?)

or  $n=0$  warps, each warp reads 32 values (corner case)

- ② write and then read in shared memory

- ③ each warp reduce 32 values, atomic add at last

?

$[B, A] \rightarrow [A]$

## 2. Thread Reduce

### \* Basic Method

- ① Each block reads 1024 consecutive columns.
- ② Each thread reduce the values it reads.

### \* Requirement for Efficiency

$$H > 1024$$

### \* Different Cases

- ① same as basic
- ② different blocks reduce the same column  
use atomic reduce for final result

### \* Corner Case / Margin Part

$$n = \lfloor a/1024 \rfloor$$

- ① threads in a block read consecutive na values
- ② each thread add up the values
- ③ atomic reduce

## 3. Warp Read Part Warp Reduce

Can't be implemented in Dace ?

# HIGH DIMENSION

\* let  $R_1 \rightarrow R_2$  to be a basic reduction with  $\geq 2$  DIMs

①  $AR_1 \rightarrow AR_2$

use  $AN$  blocks

②  $BR_1 \rightarrow R_2$

2.1) use  $BN$  blocks, reduce atomically

2.2) each thread reads  $Bn$  values and reduce first

2.3) mixture

\* if  $N_{blocks} > N_{sm}$  use 2.1

if  $N_{blocks} < N_{sm}$

use 2.2 or 2.3 to let  $N_{blocks} > N_{sm}$

\* ① the basic method depends on last two dimensions

( $B_o A_o \rightarrow A_o$  or  $A_o B_o \rightarrow B_o$ )

② discussion with  $A^*, B^*, A_o, B_o$

$$A^* = \prod_{i=1}^{N_A} A_i \quad B^* = \prod_{i=1}^{N_B} B_i$$

for 2D case, just let  $A^* = B^* = 1$

- ①  $B < 32$   $A > 1024$  : TR
- ②  $B > 32$   $A < 1024$  : WRWR
- ③  $B > 32$   $A > 1024$  : to be tested
- ④  $A^* B^*$  large  $B < 32$   $A < 1024$  : TR?
- ⑤ corner case:  $A^* B^*$  large  $B \cdot A < 32 \cdot 32$   
(optional)