# Cuda Thread Divergence

1

Minimizing thread divergence

- Warp – Set of threads that execute the same instruction at a time.

- SIMD – Single Instruction, Multiple Data
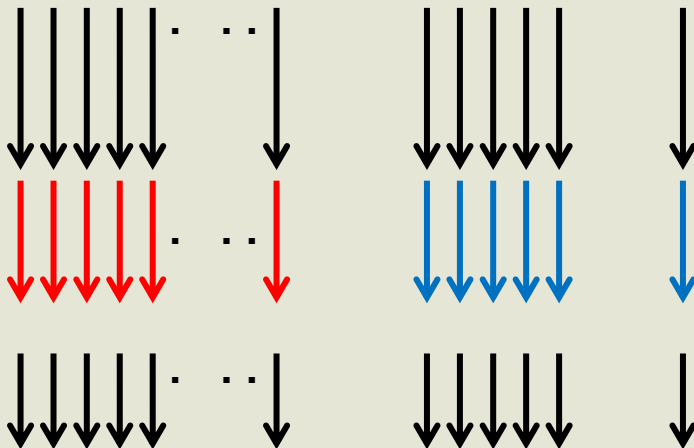
- SIMT – Singe instruction, Multiple Threads

2

```
1.   Blah;
2.   Blah;
     …
1.   If ( . . .)
2.   {
3.    //then do smth
4.   }
5.   Else
6.   {
7.   //else do smth
8.   }
9.   Blah;
10.  Blah;
```
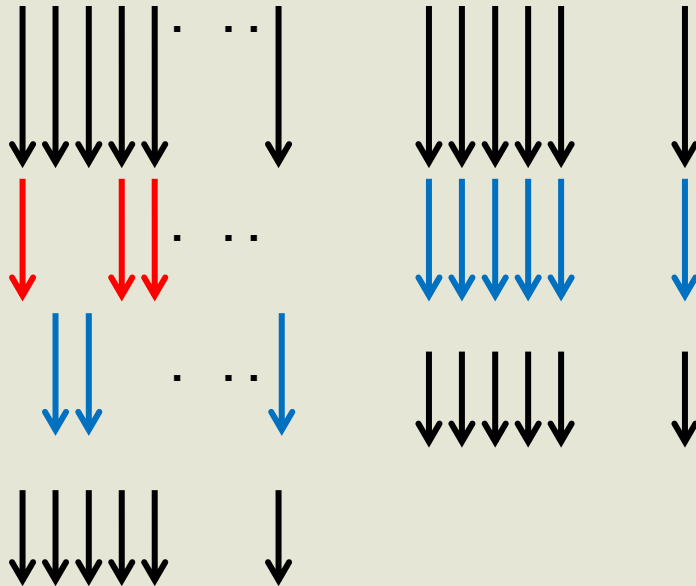
3

4

## 2-Way branch divergence



5

What is the max branch divergence penalty for a cuda thread block with 1024 threads?

_____ x slowdown

6

What is the max branch divergence penalty for a cuda thread block with 1024 threads?

# 32 x slowdown

## Max 32-way branch divergence (warp size)

7

```
1.    Switch (expr) {
2.        case 1: . . . Break;
3.        case 2: . . . Break;
4.        case 3: . . . Break;
5.        …
6.        case 32: . . . Break;
7.    }
```

8

## Survey

1. Switch (threadidx.x %32) {case 0..31}
2. Kernel <<< 1, 1024 >>>();        _____

3. Switch (threadidx.x **% 64**) {case 0..**63**}
4. Kernel <<< 1, 1024 >>>();        _____

5. Switch (**threadidx.y**) {case 0..31}
6. Kernel <<< 1, 64 x 16 >>>();        _____

7. Switch (**threadidx.y**) {case 0..31}
8. Kernel <<< 1, 16 x 16 >>>();        _____

- What will be the slowdown for each of the following expressions in switch statements

9

## Survey

1. Switch (threadidx.x % 2) {case 0..31}
2. Kernel <<< 1, 1024 >>>();        _____

3. Switch (threadidx.x **/ 32**) {case 0..**31**}
4. Kernel <<< 1, 1024 >>>();        _____

5. Switch (threadidx.x **/ 8**) {case 0..**63**}
6. Kernel <<< 1, 1024 >>>();        _____

- What will be the slowdown for each of the following expressions in switch statements

10

5

Branch Divergence in real world

- Assume a 1024 x 1024 image.

- Requiring Special handling of pixels on the boundary

11

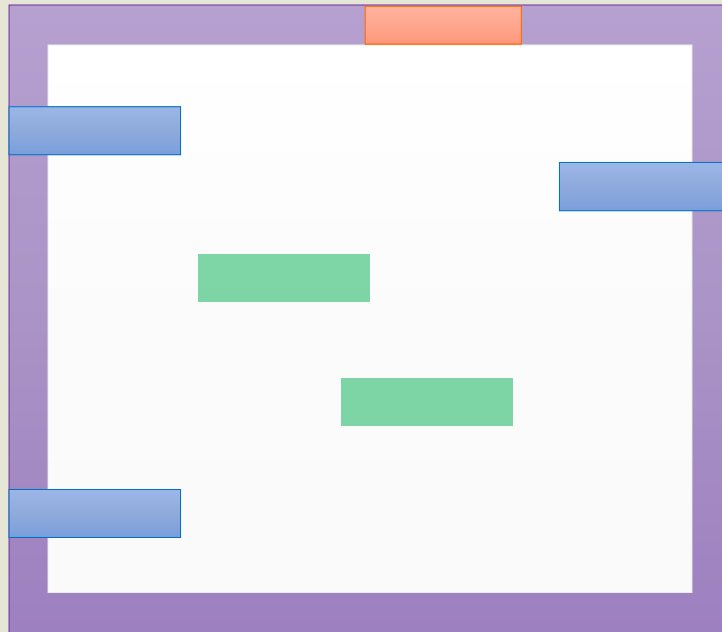Maximum branch divergence of any warp? _____ - way

```
1.    __global__
2.    If(threadIdx.x == 0 ||
         threadxIdx.x == 1024
         threadIdx.y==0 ||
         threadIdx.y==1024){
1.       //deal with boundary cond.
2.    }else {
3.       //do smth
4.    }
```

12

Branch divergence in convolution



13

Branch divergence in real life

- Be aware of branch divergence

- Don't panic if there are if statements

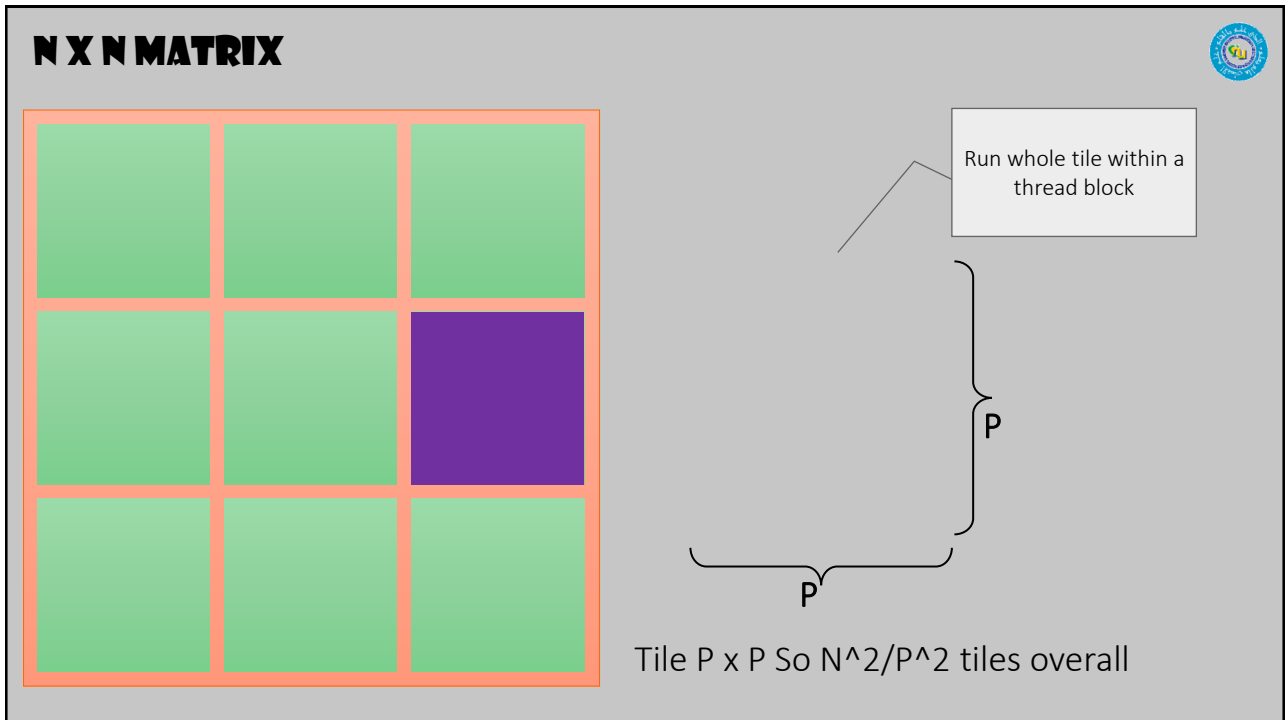- No real strategy in reducing branch divergence

14

Major guidelines

- Avoid code with too many branches

- Be aware of large imbalance in thread workloads
  - For loops with variable terminating statements.

15

# TILED MATRIX OPERATIONS

16

8

# N X N MATRIX



Run whole tile within a thread block

P

P

Tile P x P So N^2/P^2 tiles overall
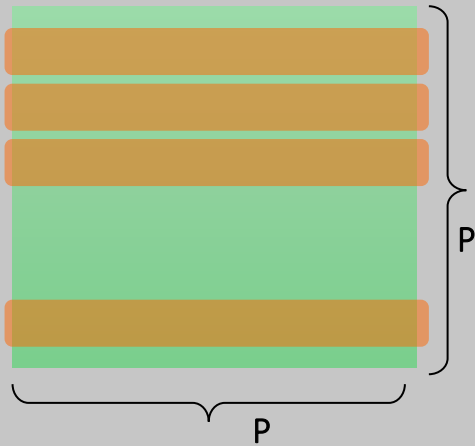
17

# P^2 THREADS TO COMPUTE 1 P X P TILE



P

P

- Good: $P^2$ parallel ops

- Bad: Must share parameters btw threads

- How many threads must get the parameters for
  - Each source element?____
  - Each dest. Element?_____

18

## P THREADS TO COMPUTE 1 P X P TILE
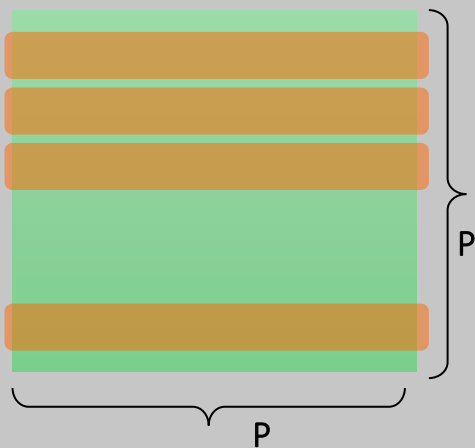
P

P

Is parallelization

increased? _____

decreased?_____

No change?_____

19

## P THREADS TO COMPUTE 1 P X P TILE

P

P

- fewer threads

- More work per thread

- Communication
  - among threads vs within a thread

20