# CS344 Introduction to Parallel Programming
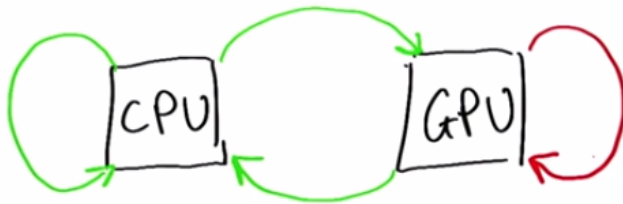
**Lesson 7.2: Dynamic Parallelism**
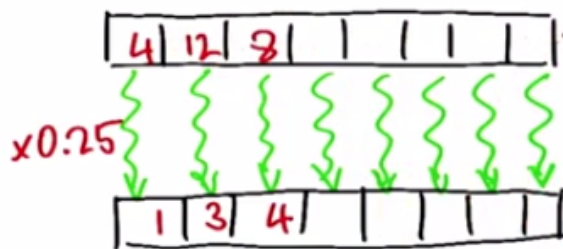
L7.2-7.1-Introduction to Dynamic Parallelism



L7.2-7.2-Bulk Parallelism



L7.2-7.3-Bulk Parallelism Quiz

## NESTED PARALLELISM



## EXAMPLE: AUDIO PROCESSING

KERNEL

① → $\{$ $\{$ $\{$ $\{$ $\{$ $\{$

MAX = 1.8V

② ÷1.8 $\{$ $\{$ $\{$ $\{$ $\{$ $\{$

NORMALIZED AUDIO

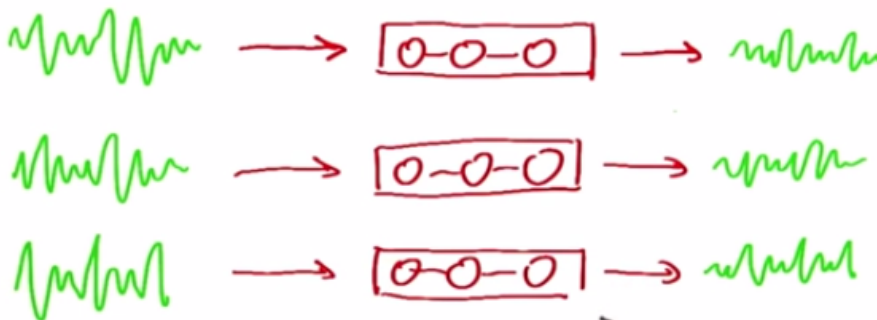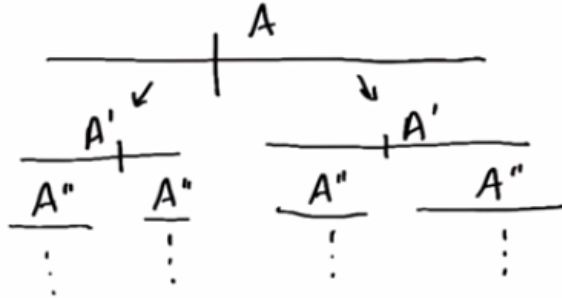L7.2-7.5-Task Parallelism

## TASK PARALLELISM

## RECURSIVE PARALLELISM

## QUIZ

WHICH TYPE OF PARALLELISM - BULK, NESTED
TASK or RECURSIVE - IS REPRESENTED
IN WHICH ALGORITHM?

|  |  |
|---|---|
|  | -MATRIX MULTIPLICATION |
|  | -FINDING FACES IN PHOTOS |
|  | -BINARY SEARCH |
|  | -CALLING A PARALLEL LIBRARY FROM A KERNEL |

## PROGRAMMING MODEL

```
__global__ void Hello() {
 printf("Hello ");
}

void main() {
 Hello<<< 1, 1 >>>();
 cudaDeviceSynchronize();
 printf("World");
}
```

```
__global__ void Hello() {
 printf("Hello ");
}

__global__ void HelloWorld() {
 Hello<<< 1, 1 >>>();
 cudaDeviceSynchronize();
 printf("World");
}
```

[CS344 Intro to Parallel Programming, Lesson 7.2]

COMPOSABILITY

PARENT   CHILD   GRANDCHILD

A   X   P

B   Y   Q

C   Z   R

THINGS TO WATCH OUT FOR

1. EVERY THREAD EXECUTES
   THE SAME PROGRAM
   - LOTS OF LAUNCHES!

QUIZ

WHAT SHOULD WE ADD TO MAKE ONLY
   THE FIRST THREAD LAUNCH THE KERNEL?

```
__global__ void launcher(){
    if([        ])
        kernel <<< 1,1 >>> ();
}
```

HINT: "threadIdx.x" gives a thread's ID

L7.2-7.12-Other Things to Watch Out For

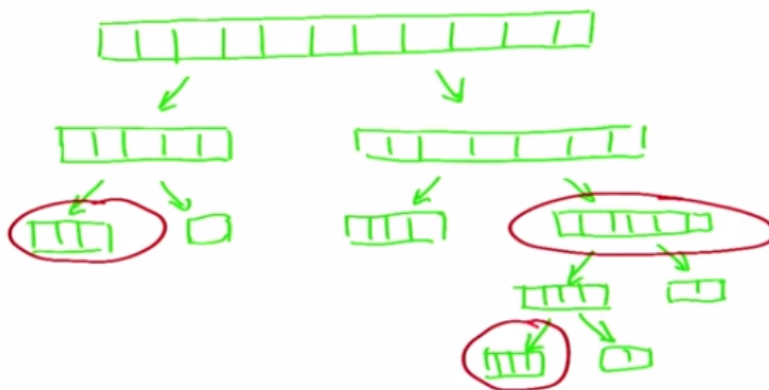## OTHER THINGS TO WATCH OUT FOR

2. EACH BLOCK EXECUTES INDEPENDENTLY
   - ALL STREAMS & EVENTS ARE PRIVATE TO THE BLOCK WHICH CREATED THEM

3. A BLOCK'S PRIVATE DATA IS PRIVATE
   - CANNOT PASS SHARED MEMORY TO CHILD KERNELS

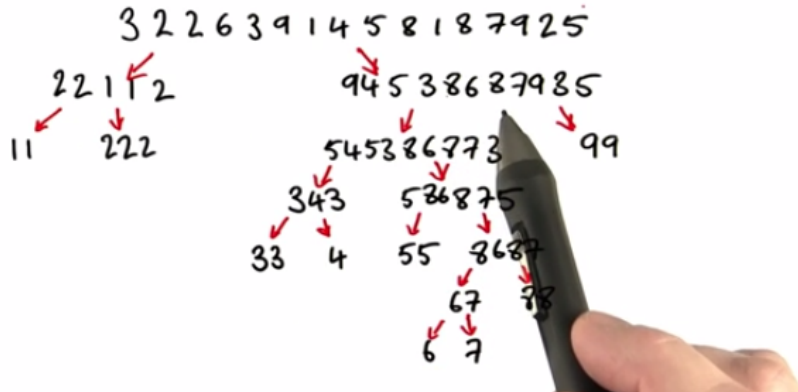L7.2-7.13-Whith Variable Cannot Be Passed to the Child Thread

## QUIZ

IN THE PROGRAM BELOW, WHICH VARIABLE MAY NOT BE PASSED TO THE CHILD KERNEL?

☐ __device__ int x[10];
☐ __shared__ float y[200];
☐ __global__ void program(){
   int *z = (int *)malloc(1000);
   launch <<<...>>> (x, y, z);
}

L7.2-7.14-Recursion and Quicksort

## RECURSION & QUICKSORT

L7.2-7.15-Dynamic Parallel Quicksort

DYNAMIC PARALLEL QUICKSORT

3 2 2 6 3 9 1 4 5 8 1 8 7 9 2 5
2 2 1 1 2              9 4 5 3 8 6 8 7 9 8 5
1 1    2 2 2           5 4 5 3 8 6 7 7 3    9 9
            3 4 3    5 3 6 8 7 5
        3 3   4    5 5   8 6 8 7
                        6 7   7 8
                      6   7

QUICKSORT EXAMPLE

```
__global__ void quicksort(int *data, int left, int right) {
    int nleft, nright;
    cudaStream_t s1, s2;

    partition(data+left, data+right, data[left], nleft, nright);

    if(left < nright){
        cudaStreamCreateWithFlags(&s1, cudaStreamNonBlocking);
        quicksort<<< ..., s1 >>>(data, left, nright);
    }
    if(nleft < right){
        cudaStreamCreateWithFlags(&s2, cudaStreamNonBlocking);
        quicksort<<< ..., s2 >>>(data, nleft, right);
    }
}
```

L7.2-7.16-Why Is Dynamic Parallel Quicksort

QUIZ

WHICH OF THE FOLLOWING REASONS
EXPLAINS WHY DYNAMIC PARALLEL QUICKSORT
IS MORE EFFICIENT?

☐ MORE EFFICIENT PARTITIONING
☑ LAUNCHING ON-THE-FLY
☐ SIMPLER CODE
☑ GREATER GPU UTILIZATION