

Parallel Computing with GPUs: Parallel

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Dr Paul Ric

<http://paulrichmond.shef.ac.uk> Add WeChat [edu_assist_pro](https://eduassistpro.github.io/) COM4521/



The
University
Of
Sheffield.

 NVIDIA

GPU
RESEARCH
CENTER

❑ Parallel Patterns Overview

❑ Reduction

❑ Scan

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

What are parallel Patterns

- ❑ Parallel patterns are high level building blocks that can be used to create algorithms
- ❑ Implementation is abstracted to give a higher level view
- ❑ Patterns describe techniques suited to parallelism
 - ❑ Allows algorithms to be built from ground up
 - ❑ Top down approach
- ❑ Consider a the simplest parallel pattern
 - ❑ Takes the input list i
 - ❑ Applies a function f
 - ❑ Writes the result list o by applying f to all members of i
 - ❑ Equivalent to a CUDA kernel where i and o are memory locations determined by `threadIdx` etc.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Gather

- ❑ Multiple inputs and single coalesced output
- ❑ Might have sequential loading or random access
 - ❑ Affect memory performance
- ❑ Differs to map due to multiple inputs

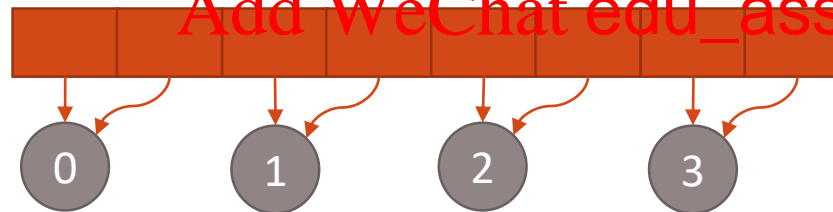
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Memory Values/Locations

ThreadId.x

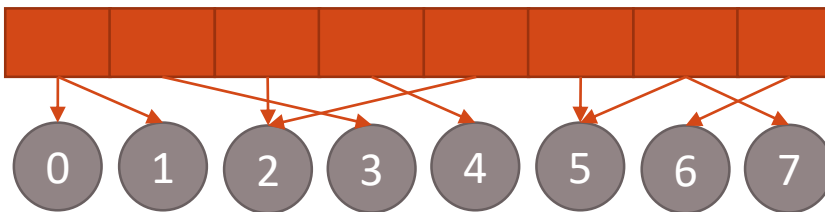


Gather operation

- ❑ Read from a number of locations

Memory Values/Locations

ThreadId.x



Gather operation

- ❑ Read from a number of locations
- ❑ Random access load

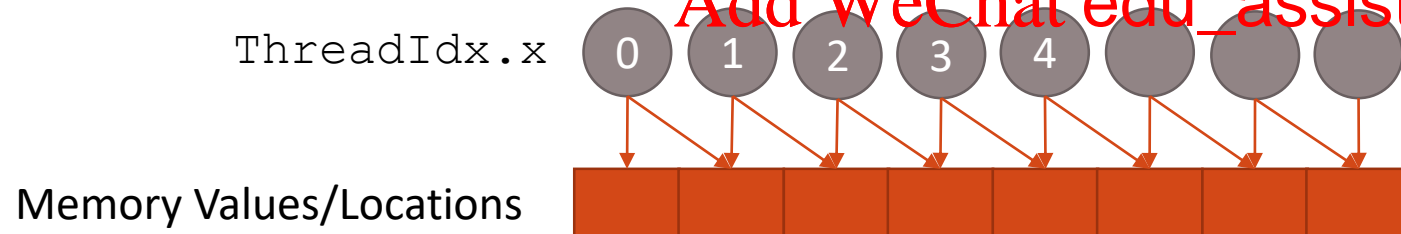
Scatter

- ❑ Reads from a single input and writes to one or many
- ❑ Can be implemented in CUDA using atomics
- ❑ Write pattern will determine performance

Assignment Project Exam Help

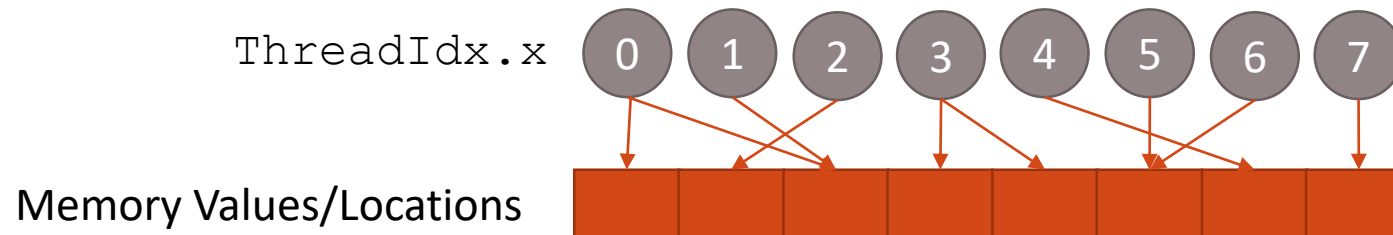
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Scatter operation

- ❑ Write to a number of locations
- ❑ Collision on write



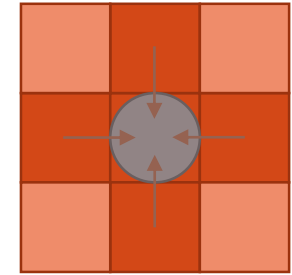
Scatter operation

- ❑ Write to a number of locations
- ❑ Random access write?

Other Parallel Patterns

❑ Stencil

- ❑ Gather a fixed pattern, usually based on locality
- ❑ See 2D shared memory examples



Stencil Gather

❑ Reduce (this lecture)

- ❑ Reduce value to a single pair
- ❑ Combined with Map (with intermediate shuffle or sort)

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

❑ Scan (this lecture)

- ❑ Compute the sum of previous value in a set

❑ Sort (*later*)

- ❑ Sort values or <value, key> pairs

❑ Parallel Patterns Overview

❑ Reduction

❑ Scan

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Reduction

❑ A reduction is where **all** elements of a set have a common *binary associative operator* (\oplus) applied to them to “reduce” the set to a single value

❑ Binary associative = order in which operations is performed on set does not matter

❑ E.g. $(1 + 2) + 3 + 4 == 1 + (2 + 3) + 4 == 10$

❑ Example operators **Assignment Project Exam Help**

❑ Most obvious example is a

❑ Other examples, Maximum <https://eduassistpro.github.io/>

❑ Serial example is trivial but how does this **allel?**

Add WeChat edu_assist_pro

```
int data[N];
int i, r;
for (int i = 0; i < N; i++){
    r = reduce(r, data[i]);
}
```

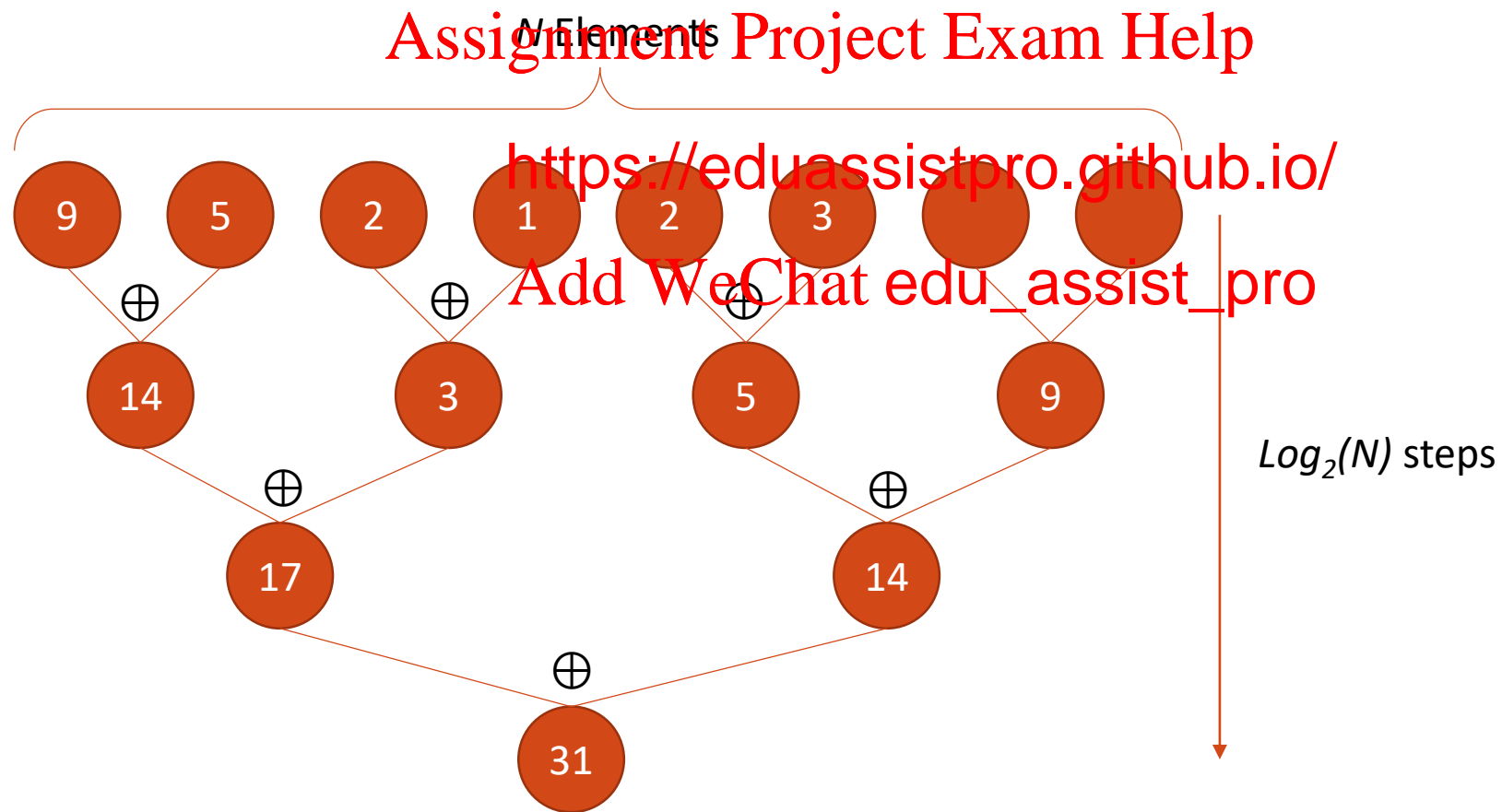
OR

```
int data[N];
int i, r;
for (int i = N-1; i >= 0; i--){
    r = reduce(r, data[i]);
}
```

```
int reduce(int r, int i){
    return r + i;
}
```

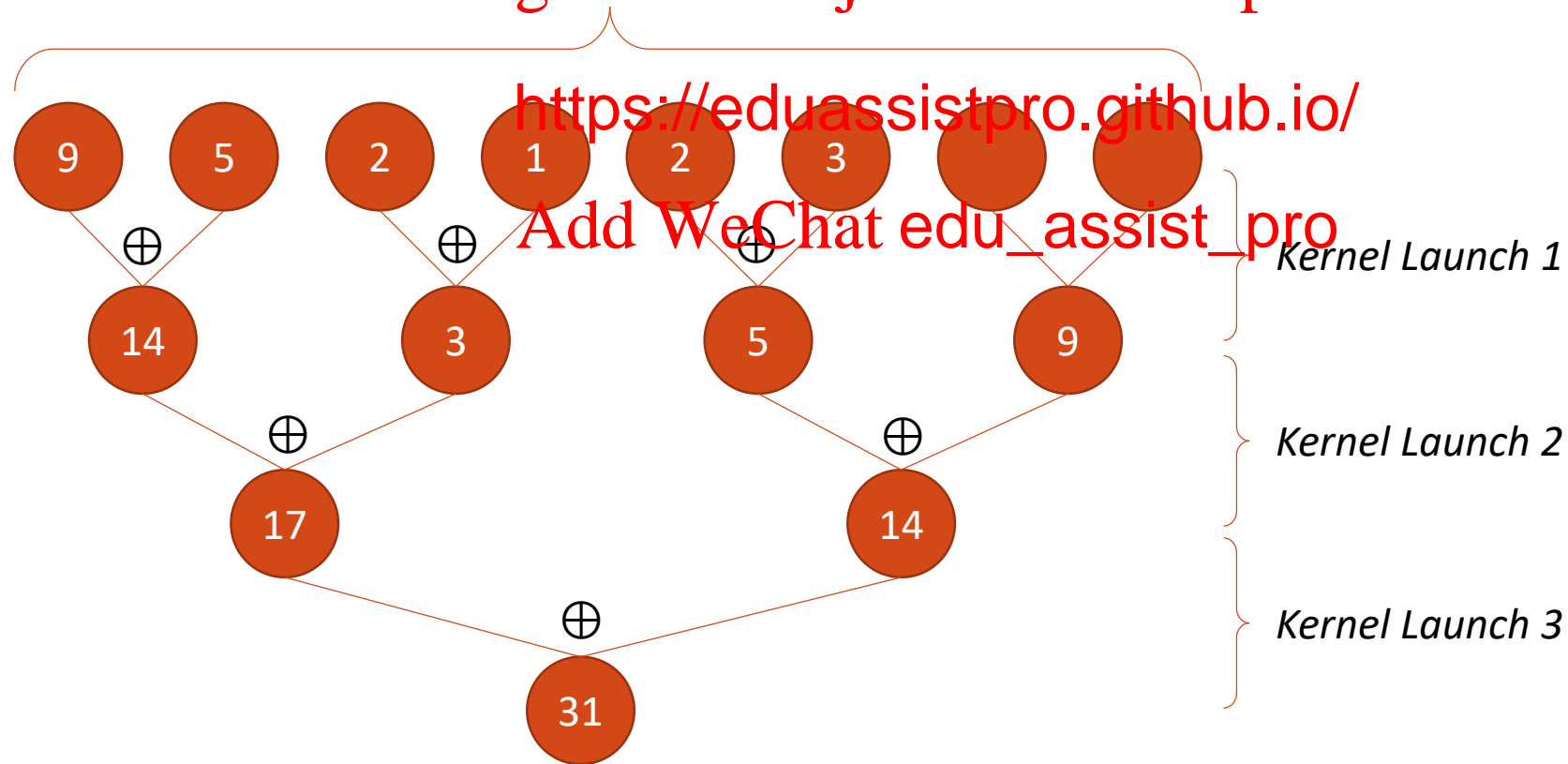

Parallel Reduction

- ❑ Order of operations does not matter so we don't have to think serially.
- ❑ A tree based approach can be used
 - ❑ At each step data is reduced by a factor of 2



Parallel Reduction in CUDA

- ❑ No global synchronisation so how do multiple blocks perform reduction?
- ❑ Split the execution into multiple stages
 - ❑ Recursive method





Recursive Reduction Problems

❑ What might be some problems with the following?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

```
__global__ void sum_reductio
    __syncthreads();
    if (i % 2 == 0) {
        results[i / 2] = sdata[threadIdx.x] + sdata[threadIdx.x+1]
    }
}
```

Block Level Reduction

- ❑ Lower launch overhead (reduction within block)
- ❑ Much better use of shared memory

```
__global__ void sum_reduction(float *input, float *block_results){
    extern __shared__ int sdata[];

    unsigned int i = blockIdx.x;
    sdata[threadIdx.x] = input[i];
    __syncthreads();

    for (unsigned int stride = blockDim.x; stride > 1; stride /= 2){
        unsigned int strided_i = threadIdx.x * 2;
        if (strided_i < blockDim.x){
            sdata[strided_i] += sdata[strided_i + stride];
        }
        __syncthreads();
    }

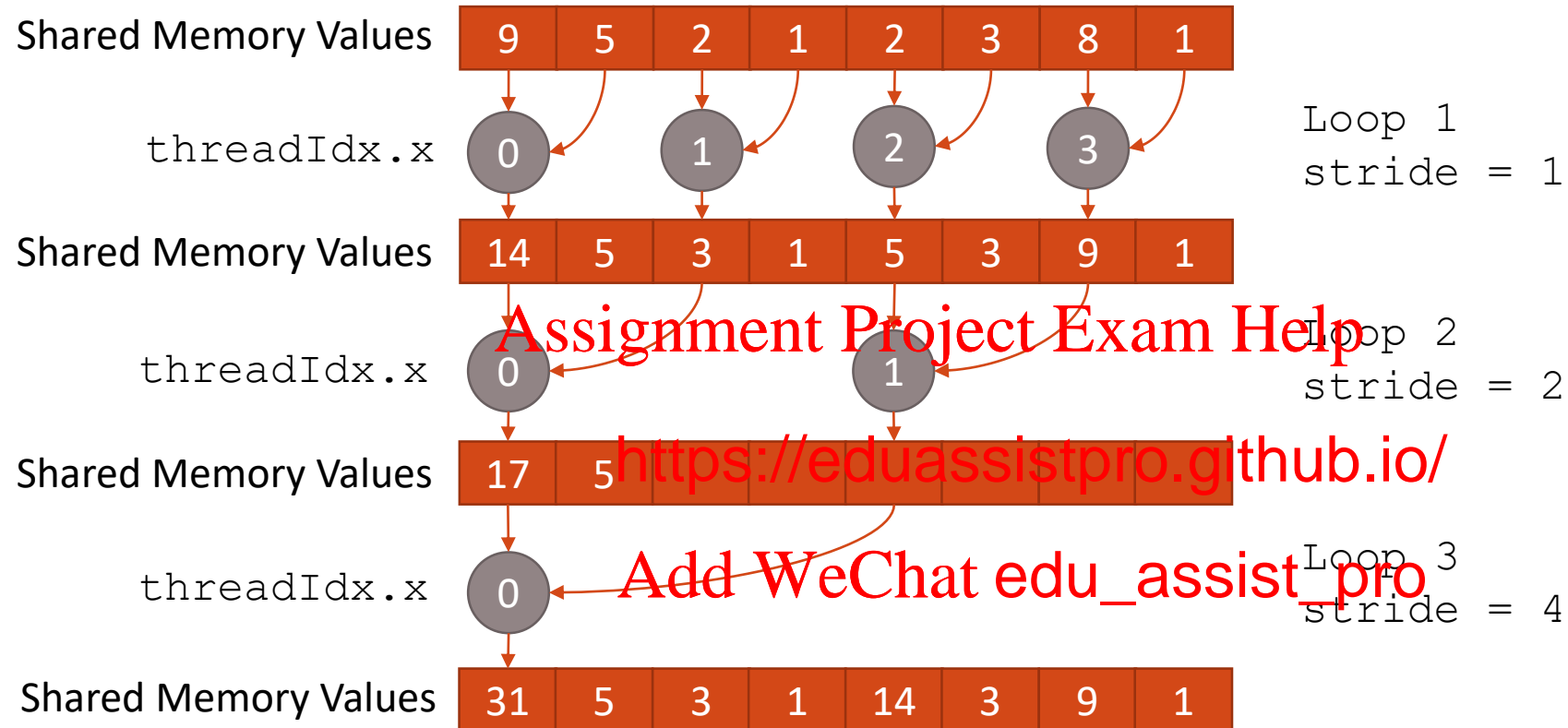
    if (threadIdx.x == 0)
        block_results[blockIdx.x] = sdata[0];
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)

Block Level Recursive Reduction



```
for (unsigned int stride = 1; stride < blockDim.x; stride*=2){  
    unsigned int strided_i = threadIdx.x * 2 * stride;  
    if (strided_i < blockDim.x){  
        sdata[strided_i] += sdata[strided_i + stride]  
    }  
    __syncthreads();  
}
```



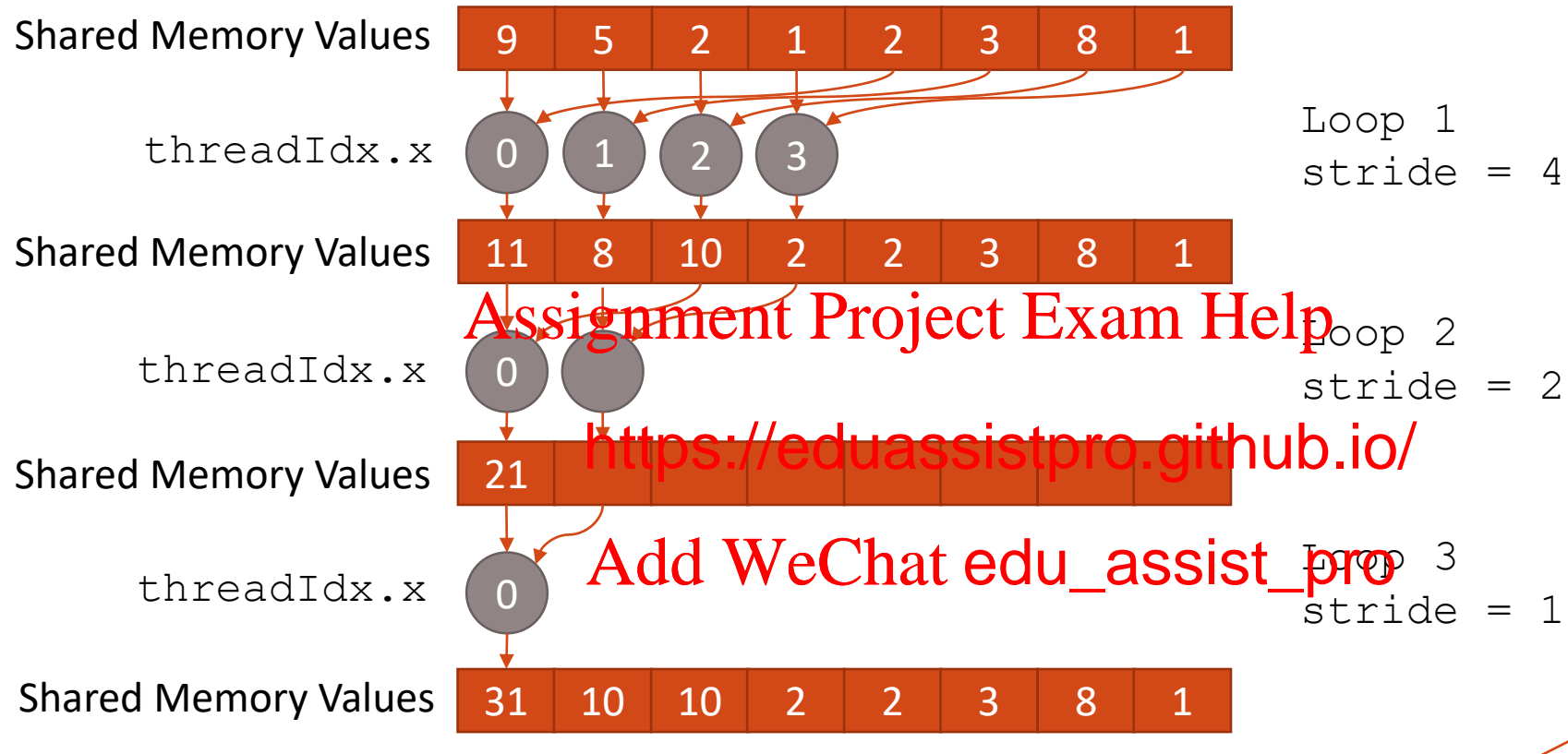
Block Level Reduction

❑ Is this shared memory access pattern bank conflict free?

Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

```
for (unsigned int stride = 1; stride < blockDim.x; stride*=2) {  
    unsigned int stride;  
    if (strided_i < blockDim.x) {  
        sdata[strided_i] = sdata[strided_i] + sdata[strided_i + stride];  
    }  
    __syncthreads();  
}
```

Block Level Reduction (Sequential Addressing)



```
for (unsigned int stride = blockDim.x/2; stride > 0; stride>>=1){  
    if (threadIdx.x < stride){  
        sdata[threadIdx.x] += sdata[threadIdx.x + stride]  
    }  
    __syncthreads();  
}
```

sm_stride	1		
loop stride	1		
threadIdx.x		index	bank
0		1	1
1		2	2
2		3	3
3		4	4
4		5	5
5		6	6
6		7	7
7		8	8
8		9	9
9		10	10
10		11	11
11		12	12
12		13	13
13		14	14
14		15	15
15		16	16
16		17	17
17		18	18
18		19	19
19		20	20
20		21	21
21		22	22
22		23	23
23		24	24
24		25	25
25		26	26
26		27	27
27		28	28
28		29	29
29		30	30
30		31	31
31		32	0
		Banks Used	32
		Max Conflicts	1

❑ Now conflict free regardless of the reduction loop stride

❑ The stride between shared memory variable accesses for threads is *always* sequential

Assignment Project Exam Help

<https://eduassistpro.github.io/>

❑ Careful types of stride discussed
Add WeChat edu_assist_pro

1. Loop stride (of algorithm)
2. SM *variable* stride (in 4 bytes)

Global Reduction Approach

- ❑ Use the recursive method

 - ❑ Our block level reduction can be applied to the result

 - ❑ At some stage it may be more effective to simply sum the final block on the CPU

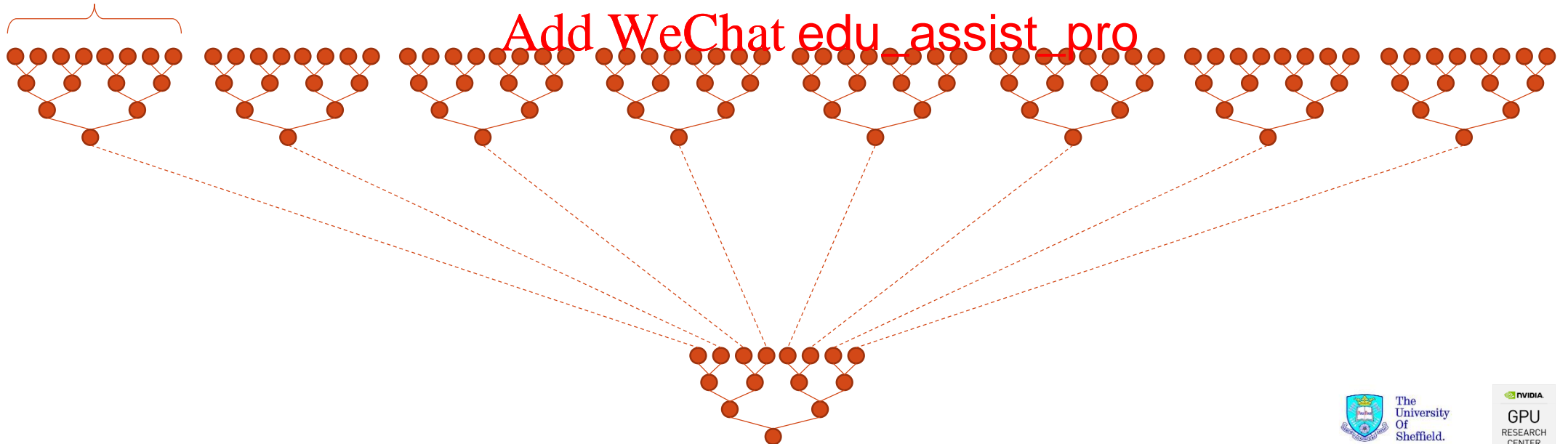
Assignment Project Exam Help

- ❑ Or use atomics on blo

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Thread block width



Global Reduction Atomics

```
__global__ void sum_reduction(float *input, float *result){
    extern __shared__ int sdata[];

    unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
    sdata[threadIdx.x] = input[i];
    __syncthreads();

    for (unsigned int stride = blockDim.x/2; stride >= 2; stride >>= 2){
        if (threadIdx.x < stride)
            sdata[threadIdx.x] += sdata[threadIdx.x + stride];
        __syncthreads();
    }

    if (threadIdx.x == 0)
        atomicAdd(result, sdata[0]);
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Further Optimisation?

❑ Can we improve our technique further?

```
__global__ void sum_reduction(float *input, float *result){
    extern __shared__ int sdata[];

    unsigned int i = blockIdx.x * blockDim.x + threadIdx.x;
    sdata[threadIdx.x] = input[i];
    __syncthreads();

    for (unsigned int stride = blockDim.x/2; stride >= 1; stride >>= 2){
        if (threadIdx.x < stride){
            sdata[threadIdx.x] += sdata[threadIdx.x + stride];
        }
        __syncthreads();
    }

    if (threadIdx.x == 0)
        atomicAdd(result, sdata[0]);
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

❑ Parallel Patterns Overview

❑ Reduction

❑ Scan

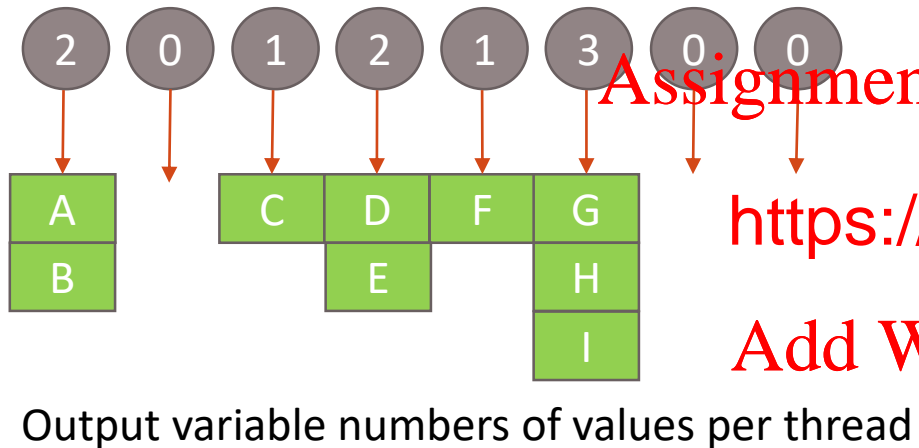
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

What is scan?

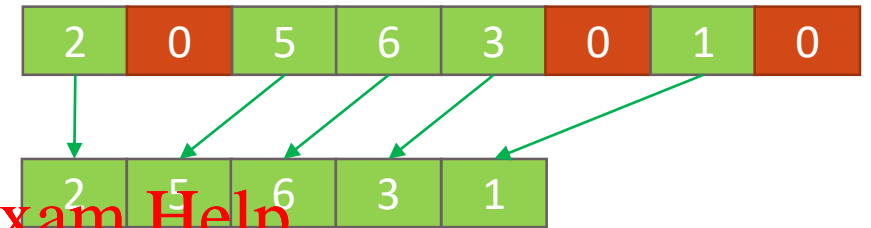
❑ Consider the following ...



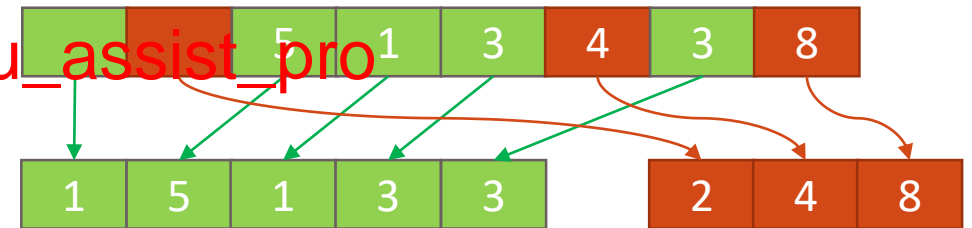
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



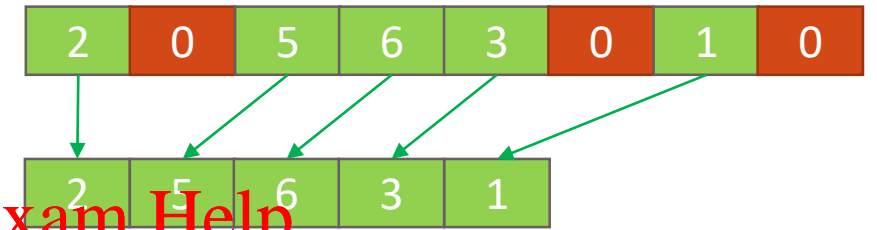
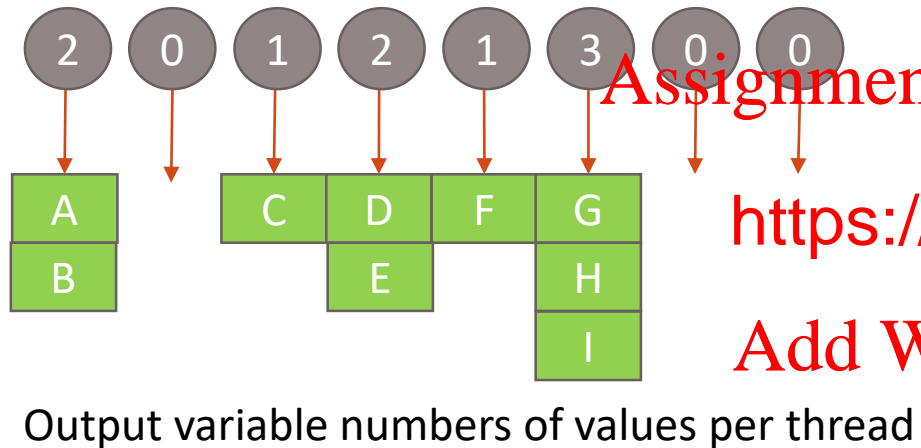
ve empty elements from array (compact)



Split elements from array based on condition (split)

What is scan?

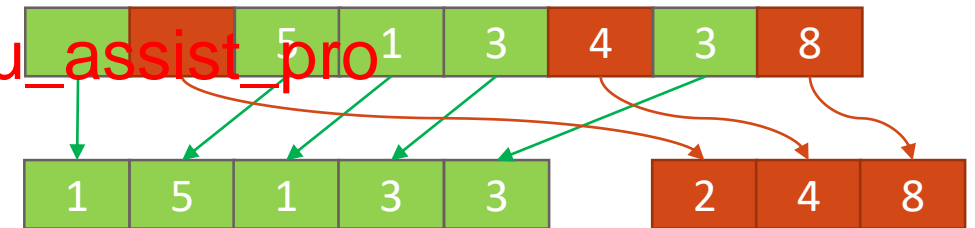
❑ Consider the following ...



Remove empty elements from array (compact)

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Split elements from array based on condition (split)

❑ Each has the same problem

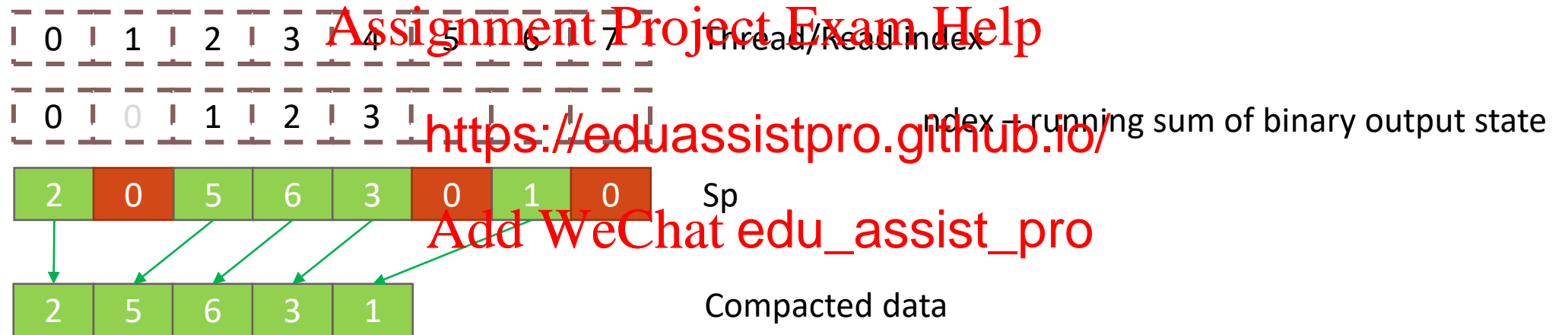
❑ Not even considered for sequential programs!

❑ Where to write output in parallel?

Parallel Prefix Sum (scan)

❑ Where to write output in parallel?

❑ Each threads needs to know the output location(s) it can write to avoid conflicts.



❑ The solution is a parallel prefix sum (or scan)

❑ Given the inputs $A = [a_0, a_1, \dots, a_{n-1}]$ and binary associate operator \oplus

❑ $\text{Scan}(A) = [0, a_0, (a_0 \oplus a_1), \dots, (a_0 \oplus a_1 \oplus \dots \oplus a_{n-1})]$

Serial Parallel Prefix Sum Example

□ E.g. Given the input and the addition operator

□ $A = [2, 6, 2, 4, 7, 2, 1, 5]$

□ $\text{Scan}(A) = [0, 2, 2+6, 2+6+2, 2+6+2+4, \dots]$

□ $\text{Scan}(A) = [0, 2, 8, 10, 14, 21, 23, 24]$

□ More generally a serial additive scan using a running sum looks like

```
int A[8] = { 2, 6, 2, 4, 7, 2, 1, 5 };
int scan_A[8];
int running_sum = 0;
for (int i = 0; i < 8; ++i)
{
    scan_A[i] = running_sum;
    running_sum += A[i];
}
```


Serial Scan for Compaction

```
int Input[8] = { 2, 0, 5, 6, 3, 0, 1, 0 };
int A[8] = { 2, 0, 5, 6, 3, 0, 1, 0 };
int scan_A[8];
int output[5]
int running_sum = 0;
```

```
for (int i = 0; i < 8; ++i){
    A[i] = Input>0;
}
```

```
for (int i = 0; i < 8; ++i){
    scan_A[i] = running_sum;
    running_sum += A[i];
}
```

```
for (int i = 0; i < 8; ++i){
    int input = Input[i];
    if (input > 0){
        int idx = scan[i];
        output[idx] = input;
    }
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

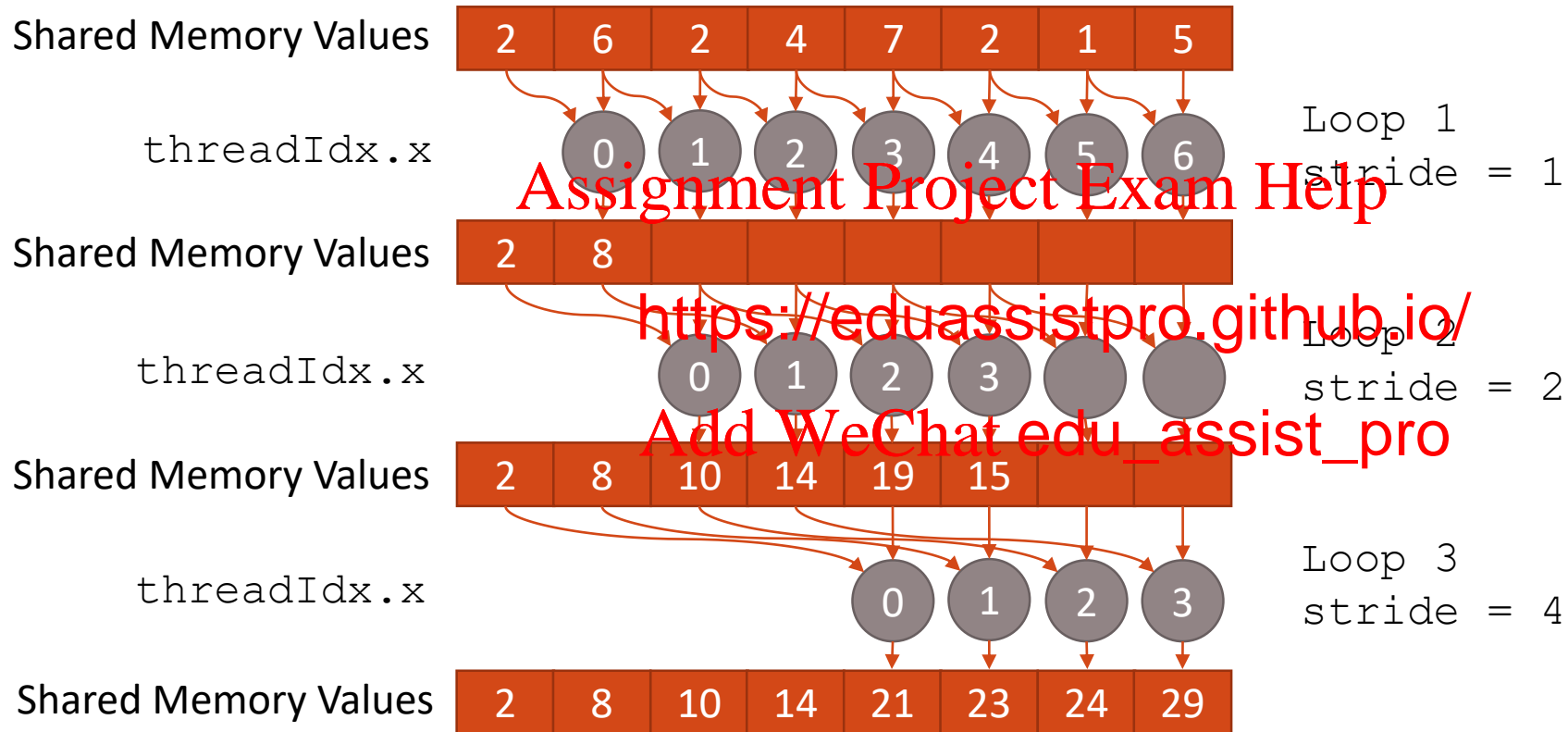
```
// generate scan input
// A = {1, 0, 1, 1, 1, 0, 1, 0}
```

```
= {0, 1, 1, 2, 3, 4, 4, 5}
```

```
// scattered write
// output = {2, 5, 6, 3, 1}
```

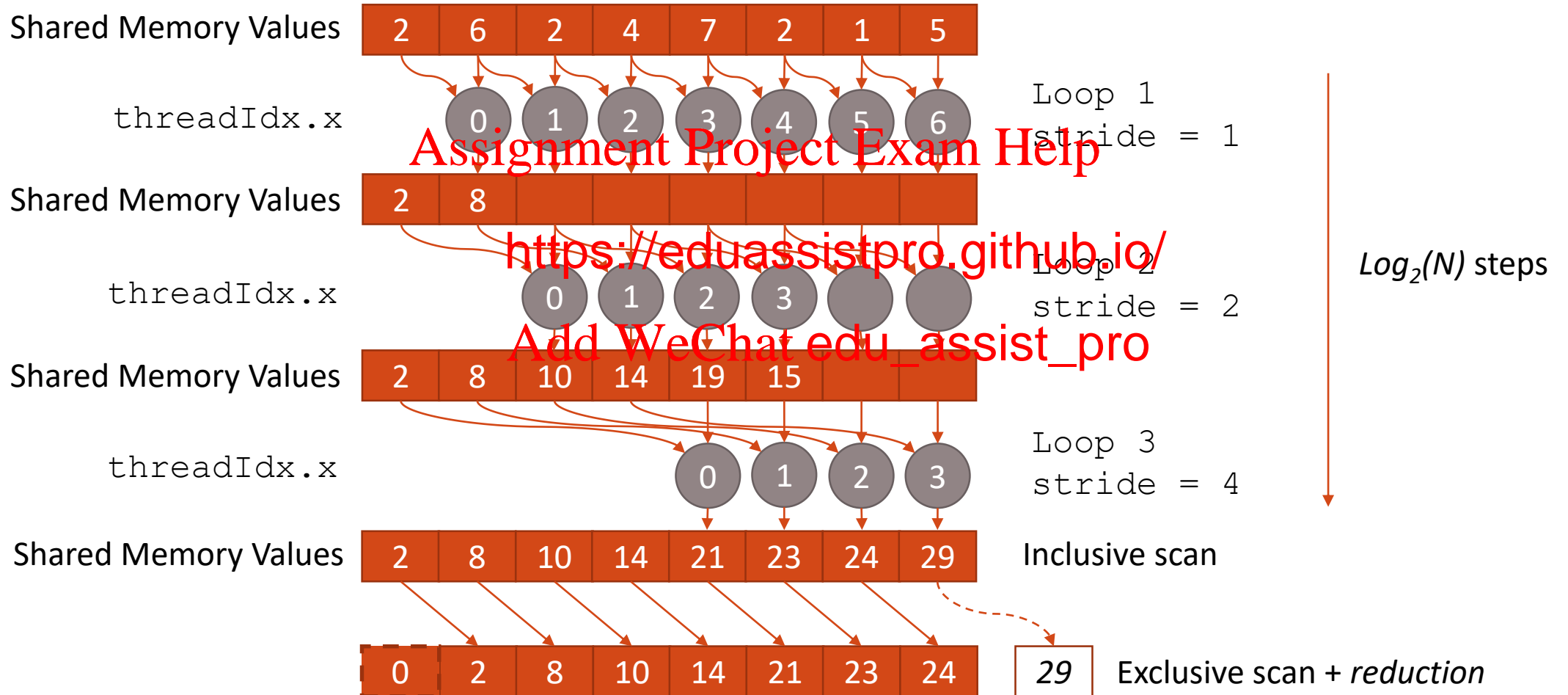
Parallel Local (Shared Memory) Scan

After $\log(N)$ loops each sum has local plus preceding $2^n - 1$ values



Inclusive Scan

Parallel Local Scan



Implementing Local Scan with Shared Memory

```
__global__ void scan(float *input) {  
    extern __shared__ float s_data[];  
    s_data[threadIdx.x] = input[threadIdx.x + blockIdx.x*blockDim.x];  
  
    for (int stride = 1; stride<blockDim.x; stride<=1) {  
        __syncthreads();  
        float s_value = (threadIdx.x > stride) ? s_data[threadIdx.x - stride] : 0;  
        __syncthreads();  
        s_data[threadIdx.x] += s_value;  
    }  
  
    //something with global results?  
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- ❑ No bank conflicts (stride of 1 between threads)
- ❑ Synchronisation required between read and write

Implementing Local Scan (at warp level)

```
__global__ void scan(float *input) {  
    __shared__ float s_data[32];  
    float val1, val2;  
  
    val1 = input[threadIdx.x + blockIdx.x*blockDim.x];  
  
    for (int s = 1; s < 32; s++) {  
        val2 = __shfl_up(val1, s);  
        if (threadIdx.x % 32 >= s)  
            val1 += val2;  
    }  
  
    //store warp level results}
```

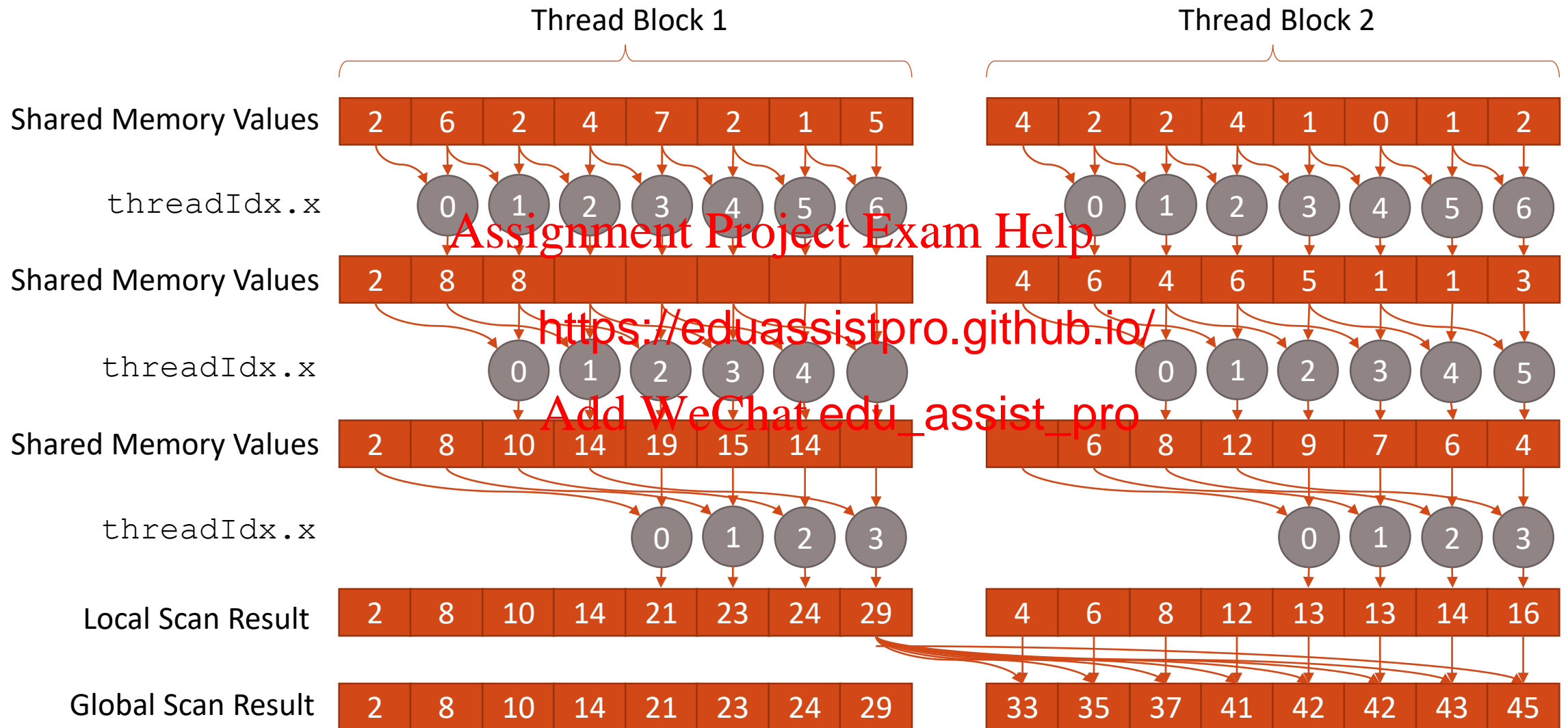
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- ❑ Exactly the same as the block level technique but at warp level
- ❑ Warp prefix sum is in `threadIdx.x % 32 == 31`
- ❑ Either use shared memory to reduce between warps
 - ❑ Or consider the following global scan approaches.

Implementing scan at Grid Level



Implementing scan at Grid Level

- ❑ Same problem as reduction when scaling to grid level
 - ❑ Each block is required to add the reduction value from proceeding blocks

- ❑ Global scan therefore requires either;
 1. Recursive scan kern
 - ❑ Additional kernel to
 2. **Atomic Increments (next slides)**
 - ❑ Increment a counter for block level results
 - ❑ Additional kernel to add sums of proceeding blocks to each value

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Global Level Scan (Atoms Part 1)

```
__device__ block_sums[BLOCK_DIM];

__global__ void scan(float *input, float *local_result) {
    extern __shared__ float s_data[];
    s_data[threadIdx.x] = input[threadIdx.x + blockIdx.x*blockDim.x];

    for (int stride = 1; stride < blockDim.x; stride <= 1) {
        __syncthreads();
        float s_value = (threadIdx.x - stride) : 0;
        __syncthreads();
        s_data[threadIdx.x] += s_value;
    }

    //store local scan result to each thread
    local_result[threadIdx.x + blockIdx.x*blockDim.x] = s_data[threadIdx.x];

    //atomic store to all proceeding block totals
    if (threadIdx.x == 0){
        for (int i=0; i<blockIdx.x; i++)
            atomicAdd(&block_sums[i], s_data[blockDim.x-1]);
    }
}
```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Global Level Scan (Atoms Part 2)

- ❑ After completion of the first kernel, block sums are all synchronised
- ❑ Use first thread in block to load block total into shared memory
- ❑ Increment local result

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

```
__device__ block_sums[BLOCK_DIM];

__global__ void scan_update(float *global_result) {
    extern __shared__ float block_total;
    int idx = threadIdx.x + blockIdx.x*blockDim.x;

    if (threadIdx.x == 0)
        block_total = block_sums[blockIdx.x];

    __syncthreads();

    global_result[idx] = local_result[idx]+block_total;
}
```

Summary

- ❑ Parallel Patterns create a bottom up model for constructing algorithms from parallel building blocks
- ❑ Reduction can be implemented recursively however re-use of data avoid costly memory movement operations
- ❑ Scan is a building block
 - ❑ Can be used for computing <https://eduassistpro.github.io/>
- ❑ Parallel patterns can be optimised read block and grid levels.
Add WeChat edu_assist_pro
- ❑ Atomics can be used in the reduction or scan value summation between blocks or warps
- ❑ Lots of potential techniques to implement and evaluate
 - ❑ Fortunately in many cases libraries and examples already exist

Acknowledgements and Further Reading

- ❑ <https://devblogs.nvidia.com/parallelforall/faster-parallel-reductions-kepler/>
 - ❑ All about application of warp shuffles to reduction
- ❑ https://stanford-cs193g-sp2010.googlecode.com/svn/trunk/lecture_6/parallel_patterns_1.ppt
 - ❑ Scan material based loosely on this lecture
- ❑ http://docs.nvidia.com/cuda/samples/6_Advanced/reduction/doc/reduction.pdf
 - ❑ Reduction material is based on this fantastic lecture by Mark Harris (NVIDIA)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro