

# OpenCL Day 3: Basic Block of Data-Parallel Programming and Applications

Cosmin Oancea and Troels Henriksen

February 4, 2019

## Course Contents

### Map and Reduce

- Types, Semantics and Properties

- Efficient Sequentialization: Work, Depth, and Brent Lemma

- Maximum Segment Sum Problem (MSSP)

- Exercise 1: OpenCL Implementation of MSSP

### Scan and Applications

- Scan: Type, Semantics, Asymptotic, Parallel Implementation

- Exercise 2: Intra-Wave Inclusive Scan Implementation

- Other Data-Parallel Operators: Scatter, Partition

- Exercise 3, Step 2: Partition2 OpenCL Implementation

### Segmented Scan and Applications

- Segmented Scan: Type, Semantics, Asymptotic, Impleme

- Application: Sparse Matrix-Vector Multiplication

- Exercise 4: Sparse Matrix-Vector Mult in OpenCL

### Optimized Partition2

# Basic Blocks of Parallel Programming: Map

**map** :  $(\alpha \rightarrow \beta, [\alpha]) \rightarrow [\beta]$  has *inherently parallel semantics*.

Applies a function to every element of the input array producing an array of equal length.

$$\begin{array}{rcl} X = & \text{map}(f, [ & a_1, \quad a_2, \quad \dots, \quad a_n \quad ] ) \\ & \downarrow \quad \quad \downarrow \quad \quad \quad \downarrow \\ X \equiv & [ & f(a_1), \quad f(a_2), \quad \dots, \quad f(a_n) \quad ] \end{array}$$

Similar:  $\text{map2}(f, [a_1, \dots, a_n], [b_1, \dots, b_n]) \equiv [f(a_1, b_1), \dots, f(a_n, b_n)]$

# Basic Blocks of Parallel Programming: Map

**map** :  $(\alpha \rightarrow \beta, [\alpha]) \rightarrow [\beta]$  has *inherently parallel semantics*.

Applies a function to every element of the input array producing an array of equal length.

$$\begin{array}{rcl} X = \text{map}(f, [ & a_1, & a_2, & \dots, & a_n & ]) \\ & \downarrow & \downarrow & & \downarrow & \\ X \equiv & [ & f(a_1), & f(a_2), & \dots, & f(a_n) & ] \end{array}$$

Similar:  $\text{map2}(f, [a_1, \dots, a_n], [b_1, \dots, b_n]) \equiv [f(a_1, b_1), \dots, f(a_n, b_n)]$

**Map Fusion:**  $\text{map}(f, \text{map}(g, A)) \equiv \text{map}(f \circ g, A)$

$$\begin{array}{rcl} A & \equiv & [ \quad a_1, \quad a_2, \quad \dots, \quad a_n \quad ] \\ X = \text{map}(g, A) & \equiv & [ \quad g(a_1), \quad g(a_2), \quad \dots, \quad g(a_n) \quad ] \\ Y = \text{map}(f, X) & \equiv & [ \quad f(g(a_1)), \quad f(g(a_2)), \quad \dots, \quad f(g(a_n)) \quad ] \\ \text{map}(f \circ g, A) & \equiv & [ \quad f(g(a_1)), \quad f(g(a_2)), \quad \dots, \quad f(g(a_n)) \quad ] \end{array}$$

Fusion is a very important optimizations; saves bandwidth!

# Basic Blocks of Parallel Programming: Reduce

**reduce** :  $((\alpha, \alpha) \rightarrow \alpha), \alpha, [\alpha]) \rightarrow \alpha$

**reduce**( $\odot, 0_\odot, [a_1, a_2, \dots, a_n]) \equiv 0_\odot \odot a_1 \odot a_2 \odot \dots \odot a_n$

where  $\odot$  is an associative binary operator (**otherwise bug!**)

$0_\odot$  is the neutral element of the monoid induced by  $\odot$

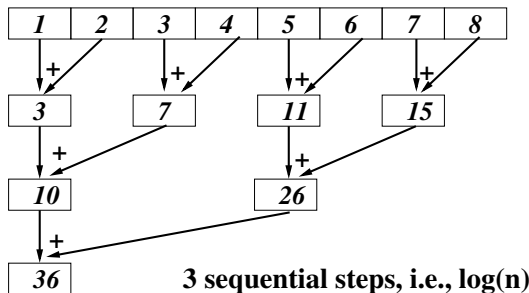
# Basic Blocks of Parallel Programming: Reduce

**reduce** :  $((\alpha, \alpha) \rightarrow \alpha), \alpha, [\alpha]) \rightarrow \alpha$

**reduce**( $\odot, 0_\odot, [a_1, a_2, \dots, a_n]) \equiv 0_\odot \odot a_1 \odot a_2 \odot \dots \odot a_n$

where  $\odot$  is an associative binary operator (**otherwise bug!**)

$0_\odot$  is the neutral element of the monoid induced by  $\odot$



Build programs by combining **map**, **reduce** and other such operators.

# Trivial Examples of Map-Reduce Programming

**Small Exercise:** write a function that receive as parameters a predicate  $p : \alpha \rightarrow \text{bool}$  and an array  $A$ , and that results in `true` if all elements satisfy  $p$  and in `false` otherwise.

Try to write (i) a divide-and-conquer and (ii) a map-reduce implementation.

# Trivial Examples of Map-Reduce Programming

**Small Exercise:** write a function that receive as parameters a predicate  $p : \alpha \rightarrow \text{bool}$  and an array  $A$ , and that results in `true` if all elements satisfy  $p$  and in `false` otherwise.

Try to write (i) a divide-and-conquer and (ii) a map-reduce implementation.

<code>all(p, [ ])</code>	<code>= True</code>	<code>all(p, x) =</code>	
<code>all(p, [x])</code>	<code>= p(x)</code>	<code>y = map(p, x)</code>	
<code>all(p, x++y)</code>	<code>= all(p, x) &amp;&amp;</code>	<code>z = reduce(&amp;&amp;, true, y)</code>	
	<code>all(p, y)</code>	<code>return z</code>	



# Trivial Examples of Map-Reduce Programming

**Small Exercise:** write a function that receive as parameters a predicate  $p : \alpha \rightarrow \text{bool}$  and an array  $A$ , and that results in `true` if all elements satisfy  $p$  and in `false` otherwise.

Try to write (i) a divide-and-conquer and (ii) a map-reduce implementation.

$\text{all}(p, []) = \text{True}$	$\text{all}(p, x) =$
$\text{all}(p, [x]) = p(x)$	$y = \text{map}(p, x)$
$\text{all}(p, x++y) = \text{all}(p, x) \ \&\&$	$z = \text{reduce}(\&\&, \text{true}, y)$
$\text{all}(p, y)$	$\text{return } z$

`++` denotes array concatenation.

A well define divide and conquer implementation requires that any split of the input array into  $x++y$  gives the same result.

(This is equivalent to the requirement that the binary operator of reduction is associative.)

Under this conditions, the two are equivalent: if you can write one, then you can derive the other (list homomorphic  $\equiv$  map-reduce).

# Asymptotic Work and Depth; Brent Lemma

Assuming an infinity of processors:

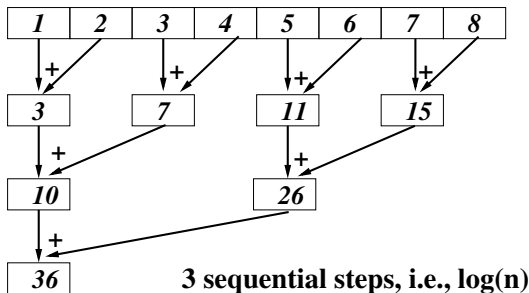
- ▶ **Work Complexity  $W(n)$** : is the total # of ops performed,
- ▶ **Depth/Step Complexity  $D(n)$** : is the # of sequential steps.
- ▶ A parallel implem is **work efficient** iff its work complexity is equal to the one of the golden sequential implem.
- ▶ Work and Depth are good high-level approximations;
- ▶ If we know the work and depth asymptotic for a program, Brent Theorem offers good complexity bounds for a PRAM.

## Theorem (Brent Theorem)

*An algorithm of depth  $D(n)$  and work  $W(n)$  can be simulated on a  $P$ -processor PRAM in time complexity  $T$  such that:*

$$\frac{W(n)}{P} \leq T \leq \frac{W(n)}{P} + D(n)$$

# Applying Brent Lemma to Map-Reduce Computation

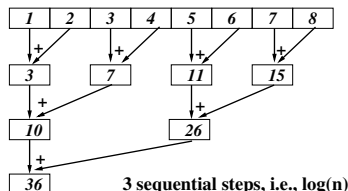


Reducing an array of length  $n$  with  $n/2$  processors requires:

- ▶ work  $W(n) = n$  and
- ▶ depth  $D(n) = \lg n$ , i.e., number of sequential steps.
- ▶ Brent Theorem states the bounds for optimal runtime  $T^{opt}$ :

$$\frac{W(n)}{p} \leq T^{opt} \leq \frac{W(n)}{p} + D(n)$$

# Applying Brent Lemma to Map-Reduce Computation



Reducing an array of length  $n$  with  $n/2$  processors requires:

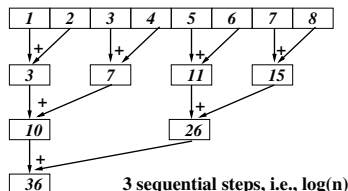
- ▶  $W(n) = n, D(n) = \lg n$ .
- ▶ Optimal runtime  $T^{opt}$  bounds:  
$$\frac{W(n)}{P} \leq T^{opt} \leq \frac{W(n)}{P} + D(n)$$

**An optimized map-reduce computation can be implemented as:**

- ▶ splits the input array into  $P$  subarrays, each containing about the same number of elements;
- ▶ perform the computation sequentially on each subarray, but in parallel across subarrays;
- ▶ reduce the  $P$  per-processor results.

**This leads to optimal runtime on  $P$  processors:**

# Applying Brent Lemma to Map-Reduce Computation



Reducing an array of length  $n$  with  $n/2$  processors requires:

- ▶  $W(n) = n, D(n) = \lg n$ .
- ▶ Optimal runtime  $T^{opt}$  bounds:  
$$\frac{W(n)}{P} \leq T^{opt} \leq \frac{W(n)}{P} + D(n)$$

**An optimized map-reduce computation can be implemented as:**

- ▶ splits the input array into  $P$  subarrays, each containing about the same number of elements;
- ▶ perform the computation sequentially on each subarray, but in parallel across subarrays;
- ▶ reduce the  $P$  per-processor results.

**This leads to optimal runtime on  $P$  processors:**  $O(\frac{n}{P} + \lg P)$

This kind of chunking is often referred to as:  
“efficient sequentialization of excess parallelism”!

# Almost Homomorphisms (Gorlatch)

“Systematic Extraction and Implementation of Divide-and-Conquer Parallelism”, Sergei Gorlatch, 1996.

**Intuition:** a non-homomorphic function  $g$  can be sometimes lifted to a homomorphic one  $f$ , by computing a baggage of *extra info*.

The initial fun obtained by projecting the homomorphic result:

$$g = \pi \circ f$$

**Maximum-Segment Sum Problem (mssp):**

Given a list of integers, find the contiguous segment of the list whose members have the largest sum among all such segments. The result is only the maximal sum (not the segment's members). For simplicity let's assume we are interested only in **positive sums**.

E.g.,  $\text{mss } [1, -2, 3, 4, -1, 5, -6, 1] = 11$   
(the corresponding segment is  $[3, 4, -1, 5]$ ).

# Maximum Segment Sum (MSSP): Preliminary Reasoning

## Maximum-Segment Sum Problem (mss):

Given a list of integers, find the contiguous segment of the list whose members have the largest sum among all such segments. The result is only the maximal sum (not the segment's members). For simplicity let's assume we are interested only in **positive sums**.

## A first straightforward/naive attempt:

```
mss [ ] = 0
mss [a] = a ↑ 0 //↑ denotes Max
mss (x ++ y) = mss(x) ??? mss(y)
```

# Maximum Segment Sum (MSSP): Preliminary Reasoning

## Maximum-Segment Sum Problem (mss):

Given a list of integers, find the contiguous segment of the list whose members have the largest sum among all such segments. The result is only the maximal sum (not the segment's members). For simplicity let's assume we are interested only in **positive sums**.

### A first straightforward/naive attempt:

```
mss [ ] = 0  
mss [a] = a ↑ 0 //↑ denotes Max  
mss (x ++ y) = mss(x) ??? mss(y)
```

x				y			
1	-2	3	4	-1	5	-6	1
mss1 = 7				mss2 = 5			

Which case is problematic?

### How to combine mss1 and mss2?

**mss1 + mss2 = 12    Incorrect!**

**max(mss1, mss2) = 7    Incorrect!**



# Maximum Segment Sum (MSSP): Preliminary Reasoning

## Maximum-Segment Sum Problem (mss):

Given a list of integers, find the contiguous segment of the list whose members have the largest sum among all such segments. The result is only the maximal sum (not the segment's members). For simplicity let's assume we are interested only in **positive sums**.

### A first straightforward/naive attempt:

```
mss [ ] = 0  
mss [a] = a ↑ 0 //↑ denotes Max  
mss (x ++ y) = mss(x) ??? mss(y)
```

x				y			
1	-2	3	4	-1	5	-6	1
mss1 = 7				mss2 = 5			

Which case is problematic?

Answer: when the segment of interest lies partly in x and partly in y!

### How to combine mss1 and mss2?

$mss1 + mss2 = 12$  **Incorrect!**

$\max(mss1, mss2) = 7$  **Incorrect!**

# Maximum Segment Sum (MSSP): A Better Reasoning

The problematic case is when the segment of interest lies partly in  $x$  and partly in  $y$ !

We need to compute extra information:

# Maximum Segment Sum (MSSP): A Better Reasoning

The problematic case is when the segment of interest lies partly in  $x$  and partly in  $y$ !

We need to compute extra information:

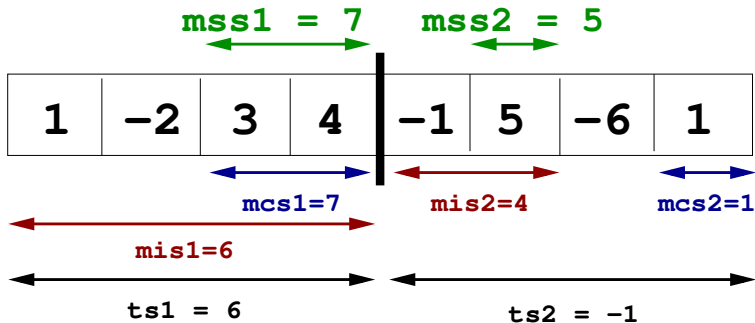
- ▶ maximum concluding segment
- ▶ maximum initial segment
- ▶ total segment sum

# Maximum Segment Sum (MSSP): A Better Reasoning

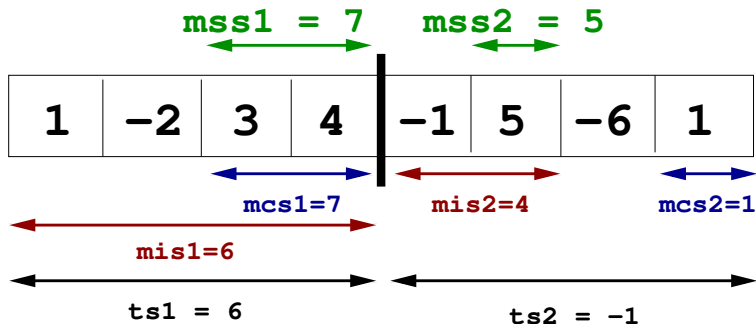
The problematic case is when the segment of interest lies partly in x and partly in y!

We need to compute extra information:

- ▶ maximum concluding segment
- ▶ maximum initial segment
- ▶ total segment sum



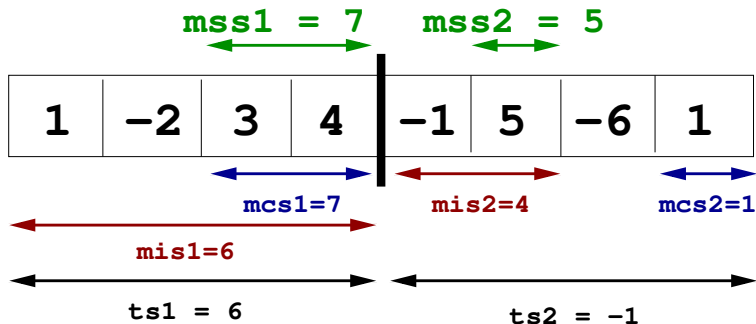
# MSSP: Deriving the Implementation



Lets compute the  $mis$ ,  $mcs$ ,  $mss$ , and  $ts$  for the result of the two concatenating segments.  $\uparrow$  denotes max.

$mis =$

# MSSP: Deriving the Implementation

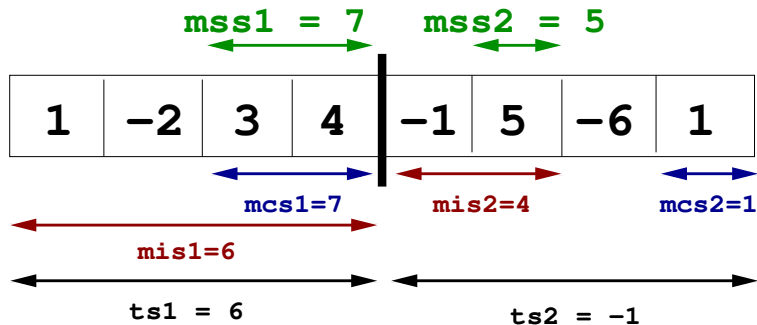


Lets compute the  $mis$ ,  $mcs$ ,  $mss$ , and  $ts$  for the result of the two concatenating segments.  $\uparrow$  denotes max.

$$mis = mis1 \uparrow (ts1 + mis2)$$

$$mcs =$$

# MSSP: Deriving the Implementation



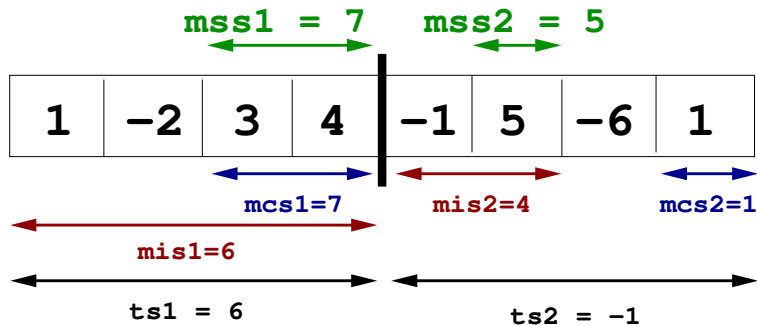
Lets compute the  $mis$ ,  $mcs$ ,  $mss$ , and  $ts$  for the result of the two concatenating segments.  $\uparrow$  denotes max.

$$mis = mis1 \uparrow (ts1 + mis2)$$

$$mcs = mcs2 \uparrow (mcs1 + ts2)$$

$$mss =$$

# MSSP: Deriving the Implementation



Lets compute the  $mis$ ,  $mcs$ ,  $mss$ , and  $ts$  for the result of the two concatenating segments.  $\uparrow$  denotes max.

$$mis = mis1 \uparrow (ts1 + mis2)$$

$$mcs = mcs2 \uparrow (mcs1 + ts2)$$

$$mss = mss1 \uparrow mss2 \uparrow (mcs1 + mis2)$$

$$ts = ts1 + ts2$$



# MSSP: Map-Reduce C-ish Pseudocode

```
//  $x \uparrow y \equiv (x \geq y) ? x : y$ 
Qtup  $\odot$  (Qtup x, Qtup y) {
    Qtup r;
    r.mss = x.mss  $\uparrow$  y.mss  $\uparrow$ 
            (x.mcs + y.mis);
    r.mis = x.mis  $\uparrow$  (x.ts + y.mis);
    r.mcs = (x.mcs + y.ts)  $\uparrow$  y.mcs;
    r.ts  = x.ts + y.ts;
    return r;
}

Qtup f(int x) {
    Qtuple r;
    int x0 = x  $\uparrow$  0;
    r.mss = x0; r.mis = x0;
    r.mcs = x0; r.ts = x;
    return r;
}

typedef struct msict {
    int mss;
    int mis;
    int mcs;
    int ts;
} Qtup;

int mss(int* xs) {
    Qtup ne, res;
    ne.mss = 0; ne.mis = 0;
    ne.mcs = 0; ne.ts = 0;

    Qtup* ys = map(f, xs);
    res = reduce( $\odot$ , ne, ys);
    return res.mss;
}
```

The baggage: 3 extra integers (misx, mcsx, tsx) and a constant number of integer operations per communication stage.

**For performance:** array ys should not be manifested in memory; fuse the map with the reduce operations.

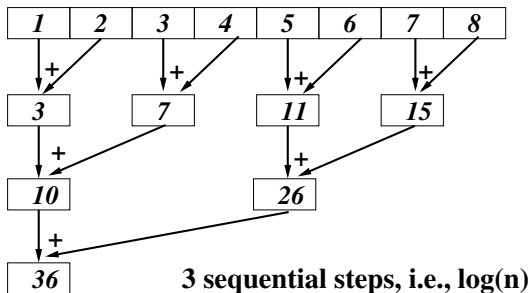
## Exercise 1: OpenCL Implementation of MSSP

- ▶ For this exercise, you will implement MSSP as an OpenCL program.
- ▶ As inspiration, we will first study four OpenCL implementations of *summing an integer array*.

$$\sum_{i < n} x[i]$$

# Binary Tree Reduction

**The idea:** each thread reads two neighbouring elements, adds them together, and writes one element. This halves the array in size. Continue until only a single element is left.



- ▶ Each level becomes a kernel invocation, with number of threads equal to half the number of array elements.
- ▶  $O(n)$  work and  $O(\log(n))$  span (optimal).
- ▶ **Why is this not efficient?**

# Improving the Tree Reduction

**The idea:** instead of shrinking the array by a factor of two for each level, shrink it by the workgroup size.

- ▶ Same asymptotic performance.
- ▶ Avoids kernels with very few threads. E.g with workgroup size 256:  $10000000 \rightarrow 39063 \rightarrow 153 \rightarrow 1$ .

# Improving the Tree Reduction

**The idea:** instead of shrinking the array by a factor of two for each level, shrink it by the workgroup size.

- ▶ Same asymptotic performance.
- ▶ Avoids kernels with very few threads. E.g with workgroup size 256:  $10000000 \rightarrow 39063 \rightarrow 153 \rightarrow 1$ .

<b>Implementation</b>	$n = 1000$	$n = 1000000$
Tree reduction	$77\mu s$	$363\mu s$
Group reduction	$17\mu s$	$179\mu s$

# Applying Brent's Lemma

**The idea:** instead of letting the thread count depend on the input size, always launch the same number of threads, and have each thread perform an efficient sequential summation of a *chunk* of the input.

- ▶ GPUs have a maximum (hardware/problem-dependent) capacity for exploiting parallelism. Beyond that limit, parallelism is at best worthless, and usually comes with overhead (e.g. excessive synchronisation).
- ▶ *A straightforward implementation of this idea only works if the operator is commutative!*

# Applying Brent's Lemma

**The idea:** instead of letting the thread count depend on the input size, always launch the same number of threads, and have each thread perform an efficient sequential summation of a *chunk* of the input.

- ▶ GPUs have a maximum (hardware/problem-dependent) capacity for exploiting parallelism. Beyond that limit, parallelism is at best worthless, and usually comes with overhead (e.g. excessive synchronisation).
- ▶ *A straightforward implementation of this idea only works if the operator is commutative!*

Implementation	$n = 1000$	$n = 1000000$
Tree reduction	$77\mu s$	$363\mu s$
Group reduction	$17\mu s$	$179\mu s$
Chunked reduction	$70\mu s$	$103\mu s$

# Using Atomics

**The idea:** GPUs have special hardware support for performing certain memory updates atomically. In OpenCL, this is exposed through *atomic operations*.

```
int atomic_add(volatile __global int *p,  
               int val)
```

- ▶ Concise parallel reduction: each thread reads an element and uses `atomic_add()` to update the same location in memory.
- ▶ **Why is this slow for large inputs?**



# Using Atomics

**The idea:** GPUs have special hardware support for performing certain memory updates atomically. In OpenCL, this is exposed through *atomic operations*.

```
int atomic_add(volatile __global int *p,  
               int val)
```

- Concise parallel reduction: each thread reads an element and uses `atomic_add()` to update the same location in memory.
- **Why is this slow for large inputs?**

Implementation	$n = 1000$	$n = 1000000$
Tree reduction	$77\mu s$	$363\mu s$
Group reduction	$17\mu s$	$179\mu s$
Chunked reduction	$70\mu s$	$103\mu s$
Atomics	$8\mu s$	$1278\mu s$

# For MSSP

Implement three versions:

1. Tree reduction
2. Group reduction
3. Chunked reduction

## Hints:

- ▶ Instead of a struct `Qtup`, you may want to use `int4` (inside the kernels) and `cl_int4` (in host code).
- ▶ The MSSP operator is *not* commutative, so you will need to use a sliding window approach for the chunked reduction.
- ▶ For the chunked reduction, you can fuse the map function into the reduction itself. For the others, a separate pass is done.

## Course Contents

### Map and Reduce

- Types, Semantics and Properties

- Efficient Sequentialization: Work, Depth, and Brent Lemma

- Maximum Segment Sum Problem (MSSP)

- Exercise 1: OpenCL Implementation of MSSP

### Scan and Applications

- Scan: Type, Semantics, Asymptotic, Parallel Implementation

- Exercise 2: Intra-Wave Inclusive Scan Implementation

- Other Data-Parallel Operators: Scatter, Partition

- Exercise 3, Step 2: Partition2 OpenCL Implementation

### Segmented Scan and Applications

- Segmented Scan: Type, Semantics, Asymptotic, Impleme

- Application: Sparse Matrix-Vector Multiplication

- Exercise 4: Sparse Matrix-Vector Mult in OpenCL

### Optimized Partition2

# Scan: A Basic Block of Data-Parallel Programming

Scan is also known as parallel prefix sum:

- ▶ computes all partial prefixes of an array;
- ▶ similar type with reduce, except that it returns an array;
- ▶ **exclusive scan**: result array starts with the neutral element;
- ▶ **inclusive scan**: starts with the first element of the input array;
- ▶ Inclusive scan is slightly more useful than exclusive scan.

# Scan: A Basic Block of Data-Parallel Programming

Scan is also known as parallel prefix sum:

- ▶ computes all partial prefixes of an array;
- ▶ similar type with reduce, except that it returns an array;
- ▶ **exclusive scan**: result array starts with the neutral element;
- ▶ **inclusive scan**: starts with the first element of the input array;
- ▶ Inclusive scan is slightly more useful than exclusive scan.

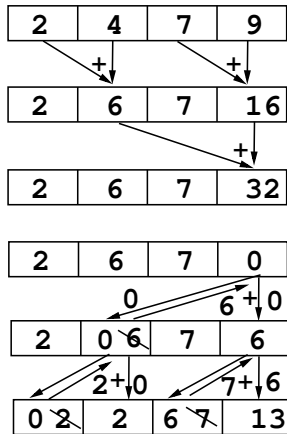
$\text{scan}^{\text{exc}}: ((\alpha, \alpha) \rightarrow \alpha), \alpha, [n]\alpha \rightarrow [n]\alpha$   
 $\text{scan}^{\text{exc}}(\odot, e, [x_1, \dots, x_n]) = [e, e \odot x_1, \dots, e \odot x_1 \odot \dots \odot x_{n-1}]$   
i.e.,  $e: \alpha, x_i: \alpha, \forall i$ , and  $\odot: (\alpha, \alpha) \rightarrow \alpha$ .

$\text{scan}^{\text{inc}}: ((\alpha, \alpha) \rightarrow \alpha), \alpha, [n]\alpha \rightarrow [n]\alpha$   
 $\text{scan}^{\text{inc}}(\odot, e, [x_1, \dots, x_n]) = [x_1, x_1 \odot x_2, \dots, x_1 \odot \dots \odot x_n]$   
i.e.,  $e: \alpha, x_i: \alpha, \forall i$ , and  $\odot: (\alpha, \alpha) \rightarrow \alpha$ .

# Parallel Exclusive Scan with Associative Operator $\oplus$

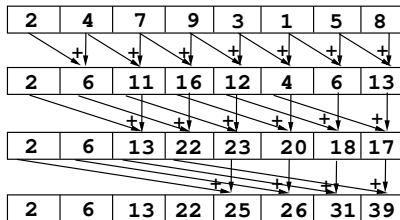
Two Steps:

- ▶ **Up-Sweep**: similar with reduction
- ▶ Root is replaced with neutral element.
- ▶ **Down-Sweep**:
  - ▶ the left child sends its value to parent and updates its value to that of parent.
  - ▶ the right-child value is given by  $\oplus$  applied to the left-child value and the (old) value of parent.
  - ▶ note that the right child is in fact the parent, i.e., in-place algorithm.



Scan's Work and Depth:  $D(n) = \Theta(\lg n)$ ,  $W(n) = \Theta(n)$

# Wavefront-Level Inclusive Scan for GPUs



Input: array A of  $n=2^k$  elements  
of type  $\alpha$

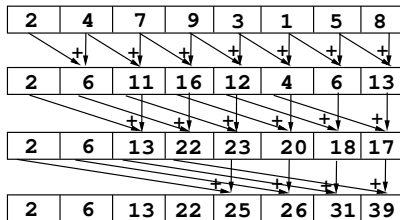
$\oplus : (\alpha, \alpha) \rightarrow \alpha$  associative

Output:  $B = [a_1, a_1 \oplus a_2, \dots, \oplus_{j=0}^{n-1} a_j]$

1. forall  $i = 0 : n-1$  do
2.      $B[i] \leftarrow A[i]$
3. endfor
4. for  $d = 0$  to  $k-1$  do
5.      $h = 2^d$
6.     forall  $i = h$  to  $n-1$  do
7.          $B[i] \leftarrow B[i-h] \oplus B[i]$
8.     endfor
9. endfor

Offers better performance because it operates in one sweep rather than two!

## Exercise 2: Intra-Wave Inclusive Scan Implementation



Input: array A of  $n=2^k$  elements  
of type  $\alpha$

$\oplus : (\alpha, \alpha) \rightarrow \alpha$  associative

Output:  $B = [a_1, a_1 \oplus a_2, \dots, \oplus_{j=0}^{n-1} a_j]$

```

1. forall i = 0 : n-1 do
2.   B[i] ← A[i]
3. endfor
4. for d = 0 to k-1 do
5.   h = 2d
6.   forall i = h to n-1 do
7.     B[i] ← B[i-h] ⊕ B[i]
8.   endfor
9. endfor
    
```

- Open file "scanapps.cl", and implement function named "incScanWave" (follow the instructions)
- Your  $n = \text{WAVE}$  and  $k = \lg \text{WAVE}$ ; Ignore the init. loop;
- Unroll the for d loop (#pragma unroll);
- loop forall i = h to n-1 is implicit,
- it should be replace by a condition if  $(i \geq h) \{ \dots \}$ ,
- except that i is not exactly the thread id.
- Remember, you want to scan each wave, independently!



# OpenCL Scan Implementation

## BLACKBOARD!

- ▶ CPU stub is in `Day3-Exercises/ScanApps/scan.h`
- ▶ Hierarchical design:
  - 1 sequentially scan `ELEMS_PER_THREAD` by each thread,
  - 2 publish in local memory
  - 3 scan at wavefront level
  - 4 “update” the scan at workgroup level
  - 5 Gather the last elements in the scan groups into a separate array, and scan that!
  - 5 Virtualize steps [1-4] so that the scan at step 5 fits into one workgroup;
  - 6 add element  $i-1$  of the result of step 5 to each element of workgroup  $i$  resulted from step 4!
- ▶ This requires 2 reads and 2 writes from/to global memory.  
Can you do it better?

# Scatter: Parallel Write Operator

**Scatter** updates in parallel and in place an input array with a set of values at specified indices:

$\text{scatter} : ([m]_{\alpha}, [n]_{\text{int}}, [n]_{\alpha}) \rightarrow [m]_{\alpha}$

A (data vector) = [b0, b1, b2, b3]

I (index vector) = [2, 4, 1, -1]

X (input array) = [a0, a1, a2, a3, a4, a5]

$\text{scatter } X \mid A$  = [a0, b2, b0, a3, b1, a5]

## Scatter: Parallel Write Operator

**Scatter** updates in parallel and in place an input array with a set of values at specified indices:

$\text{scatter} : ([m]_{\alpha}, [n]_{\text{int}}, [n]_{\alpha}) \rightarrow [m]_{\alpha}$

A (data vector) = [b0, b1, b2, b3]

I (index vector) = [2, 4, 1, -1]

X (input array) = [a0, a1, a2, a3, a4, a5]

$\text{scatter } X \mid A = [a0, b2, b0, a3, b1, a5]$

**scatter** has  $D(n) = \Theta(1)$  and  $W(n) = \Theta(n)$ ,  
i.e., requires  $n$  update operations ( $n$  is the size of  $I$  or  $A$ , not of  $X$ !).

**replicate(n, v)** creates an array of length  $n$  filled with element  $v$ !

## Exercise 3: Partition2 Operator

### Type and Semantics of Partition2:

`partition2` :  $(\alpha \rightarrow \text{Bool}, [n]\alpha) \rightarrow (\text{int}, [n]\alpha)$

Partition2 receives as input a predicate and an array and results in:

- ▶ an integer denoting the number of elements that succeed under predicate, tupled with
- ▶ a new array, having the same elements as the input array, but reordered such as the elements that succeed under the predicate occur before the others.
- ▶ The partial order of the elements that succeed/fail should be respected.

```
bool even(int v) { return (bool)(1 -(v&1)); }  
partition2(even, [5, 4, 2, 10, 3, 7, 8] ) should result in  
          (4, [4, 2, 10, 8, 5, 3, 7])
```

## Exercise 3, Step 1: Partition2 High-Level Implem

```
bool even(int v) { return (bool)(1 -(v&1)); }  
partition2(even, [5, 4, 2, 10, 3, 7, 8] ) should result in  
          (4, [4, 2, 10, 8, 5, 3, 7])
```

Step 1: implement partition2 based on map, scan and scatter (?)

```
partition2(X: [n]i32): (i32, [n]i32) =  
  cs = map(even, X);  
  tfs= map (\c->if c then 1 else 0  
           , cs);
```

## Exercise 3, Step 1: Partition2 High-Level Implem

```
bool even(int v) { return (bool)(1 -(v&1)); }  
partition2(even, [5, 4, 2, 10, 3, 7, 8] ) should result in  
      (4, [4, 2, 10, 8, 5, 3, 7])
```

Step 1: implement partition2 based on map, scan and scatter (?)

```
partition2(X: [n]i32): (i32, [n]i32) =  
  cs = map(even, X);  
  tfs= map (\c->if c then 1 else 0  
           , cs);  
  isT= scan( (+), 0, tfs );  
  i  = isT[n-1];
```

```
  ffs= map (\c->if c then 0 else 1  
           , cs );  
  tmps = scan( (+), 0, ffs );  
  isF= map (\t -> t+i, tmps);
```

```
  inds=map3(\(c,iT,iF) ->  
            if c then iT-1 else iF-1  
            , cs, isT, isF);  
  res = scatter(scratch n, inds, X);  
  return (i, res);
```

```
X    = [5, 4, 2, 10,3, 7, 8]  
n    = 7  
cs   = [F, T, T, T, F, F, T]  
tfs  = [0, 1, 1, 1, 0, 0, 1]
```

```
isT  = [0, 1, 2, 3, 3, 3, 4]  
i    = 4
```

```
ffs  = [1, 0, 0, 0, 1, 1, 0]  
tmps = [1, 1, 1, 1, 2, 3, 3]  
isF  = [5, 5, 5, 5, 6, 7, 7]
```

```
inds = [4, 0, 1, 2, 5, 6, 3]  
X    = [5, 4, 2, 10,3, 7, 8]  
Res  = [_ , _ , _ , _ , _ , _ , _]  
-----
```

```
Res  = [4, 2, 10, 8,5, 3, 7]
```

## Exercise 3: Implement Partition2 in OpenCL

- ▶ CPU stub is  
Day3-Exercises/ScanApps/partition2.h, function  
runPartition, but nothing to do there!
- (1) In Day3-Exercises/ScanApps/scanapps.cl,  
implement from scratch kernel mapPredPartKer which is  
supposed to compute tfs and ffs.
- (2) the two scans are already implemented for you!
- (3) Then implement kernel scatterPartKer which is  
supposed to compute the last part! Note that isT and isF  
are the results of scan, i.e., isF is actually tmps in slides!  
It is Ok to redundantly compute the condition again.

# What is a Segmented Scan?



# What is a Segmented Scan?

Assume an **irregular** matrix (two-dimensional array), i.e., rows have different number of elements.

**Segmented Scan Intuition:** the operation that scans each row of an irregular matrix and returns the result.

**Flat Representation of a 2D irregular array:** flag + value arrays:

# What is a Segmented Scan?

Assume an **irregular** matrix (two-dimensional array), i.e., rows have different number of elements.

**Segmented Scan Intuition:** the operation that scans each row of an irregular matrix and returns the result.

**Flat Representation of a 2D irregular array:** flag + value arrays:

[ [1, 2, 3]	Flag Array:
, [4, 5, 6, 7, 8]	[1, 0, 0, 1, 0, 0, 0, 0, 1, 0 ]
, [9, 10]	Value Array:
]	[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

The **flag array** marks with one (or something different than zero) the start of a row/segment, the other elements are zero.

The **value array**: flat sequence of elements.

$\text{sgmScan}^{inc/exc}: ((\alpha, \alpha) \rightarrow \alpha), \alpha, [n]\text{int}, [n]\alpha \rightarrow [n]\alpha$

$\text{sgmScan}^{inc}((+), 0, [1, 0, 0, 1, 0, 0, 0, 0, 1, 0]$   
 $\quad, [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]) \equiv$

# What is a Segmented Scan?

Assume an **irregular** matrix (two-dimensional array), i.e., rows have different number of elements.

**Segmented Scan Intuition:** the operation that scans each row of an irregular matrix and returns the result.

**Flat Representation of a 2D irregular array:** flag + value arrays:

[ [1, 2, 3]	Flag Array:
, [4, 5, 6, 7, 8]	[1, 0, 0, 1, 0, 0, 0, 0, 1, 0 ]
, [9, 10]	Value Array:
]	[1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

The **flag array** marks with one (or something different than zero) the start of a row/segment, the other elements are zero.

The **value array**: flat sequence of elements.

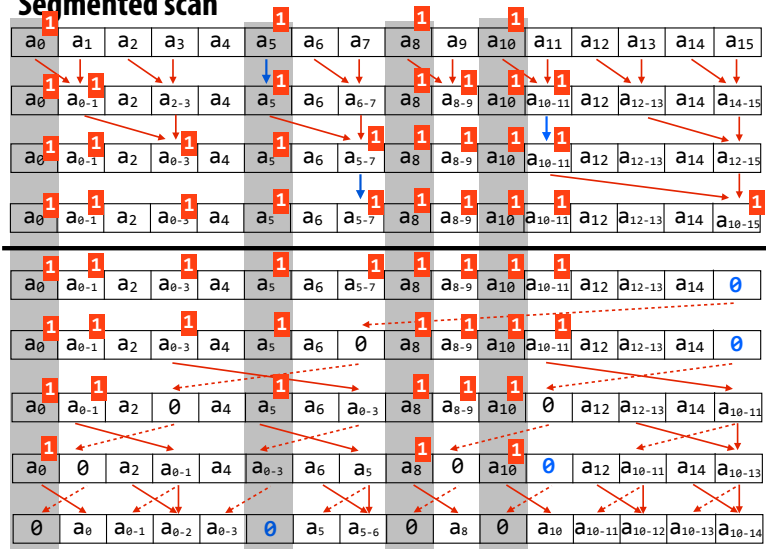
$\text{sgmScan}^{inc/exc}: ((\alpha, \alpha) \rightarrow \alpha), \alpha, [n]\text{int}, [n]\alpha \rightarrow [n]\alpha$

$\text{sgmScan}^{inc}((+), 0, [1, 0, 0, 1, 0, 0, 0, 0, 1, 0]$   
 $\quad, [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]) \equiv$   
 $\quad [1, 3, 6, 4, 9, 15, 22, 30, 9, 19]$

# Intuition: How Does the Implementation Looks Like?

Slide taken from CMU 15-418: Parallel Computer Architecture and Programming (Spring 2012). Segmented Exclusive Scan:

## Segmented scan



# Segmented Exclusive Scan Algorithm And Complexity

- ▶ While there will be more branches, the asymptotic is the same as the one for reduce and scan:
- ▶  $D(n) = \Theta(\lg n)$ ,  
 $W(n) = \Theta(n)$ !
- ▶ That's the good news!
- ▶ The bad news:

# Segmented Exclusive Scan Algorithm And Complexity

- While there will be more branches, the asymptotic is the same as the one for reduce and scan:

- $D(n) = \Theta(\lg n)$ ,  
 $W(n) = \Theta(n)$ !

- That's the good news!

- The bad news:

Input: flag array  $F$  of  $n=2^k$  of ints  
data array  $A$  of  $n=2^k$  elems of type  $T$   
 $\oplus :: T \times T \rightarrow T$  associative

Output:  $B$  = segmented scan of 2-dim array  $A$

```
1.  FORALL i = 0 to n-1 do B[i] ← A[i] ENDDO
2.  FOR d = 0 to k-1 DO // up-sweep
3.    FORALL i = 0 to n-1 by  $2^{d+1}$  DO
4.      IF  $F[i+2^{d+1}-1] == 0$  THEN
5.         $B[i+2^{d+1}-1] \leftarrow B[i+2^d-1] \oplus B[i+2^{d+1}-1]$ 
6.      ENDIF
7.       $F[i+2^{d+1}-1] \leftarrow F[i+2^d-1] \mid F[i+2^{d+1}-1]$ 
8.    ENDDO ENDDO
9.   $B[n-1] \leftarrow 0$ 
10. FOR d = k-1 downto 0 DO // down-sweep
11.   FORALL i = 0 to n-1 by  $2^{d+1}$  DO
12.     tmp ←  $B[i+2^d-1]$ 
13.     IF Foriginal $[i+2^d] \neq 0$  THEN
14.        $B[i+2^{d+1}-1] \leftarrow 0$ 
15.     ELSE IF  $F[i+2^d-1] \neq 0$  THEN
16.        $B[i+2^{d+1}-1] \leftarrow tmp$ 
17.     ELSE  $B[i+2^{d+1}-1] \leftarrow tmp \oplus B[i+2^{d+1}-1]$ 
18.     ENDIF
19.      $F[i+2^{d+1}-1] \leftarrow 0$ 
20.   ENDDO ENDDO
```

**Segmented Inclusive Scan is a Sort of Scan!**

# Segmented Inclusive Scan is a Sort of Scan!

Can be implemented as a scan with a modified binary-associative operator that operates on boolean-value (or int-value) tuples.

We need to define `zip` first, which performs SoA to AoS transform:

`zip: ( [n] $\alpha$ , [n] $\beta$  )  $\rightarrow$  [n]( $\alpha, \beta$ )`

`zip([x1, ..., xn], [y1, ..., yn]) = [(x1, y1), ..., (xn, yn)]`



# Segmented Inclusive Scan is a Sort of Scan!

Can be implemented as a scan with a modified binary-associative operator that operates on boolean-value (or int-value) tuples.

We need to define `zip` first, which performs SoA to AoS transform:

```
zip: ( [n] $\alpha$ , [n] $\beta$  )  $\rightarrow$  [n]( $\alpha, \beta$ )  
zip([x1, ..., xn], [y1, ..., yn]) = [(x1, y1), ..., (xn, yn)]
```

```
typedef Tuple {  
     $\alpha$     val;  
    bool flg;  
} Tup;
```

```
Tup  $\odot^{sgm}$ (Tup x, Tup y) {  
    Tup r;  
    r.flg = x.flg || y.flg;  
    r.val = y.flg ? y.val : x.val  $\odot$  y.val;  
    return r;  
}
```

```
 $\alpha^*$  sgmScaninc ( $\odot$ : ( $\alpha, \alpha$ )  $\rightarrow$   $\alpha$ ,  $\alpha$  ne,  $\alpha^*$  flags,  $\alpha^*$  vals) {  
    Tup nes; nes.flg = false; nes.val = ne;  
    Tup* X = zip(flags, vals);  
    return scaninc( $\odot^{sgm}$ , nes, X)  
}
```

# Sparse Matrix-Vector Multiplication: Pseudocode

```
smvMul( mat_cols: [N]int, mat_vals: [N]real
        , vct: [vct_len], shape: [num_rows]int ) : [num_rows]real =
    real res_vct[num_rows];
    int offset = 0;
    for(int i=0; i < num_rows; i++) {
        real res = 0;
        for(int j=offset; j < offset+shape[i]; j++) {
            int vct_ind = mat_cols[j];
            int mat_elm = mat_vals[j];
            res += mat_elm * vct[vct_ind];
        }
        res_vct[i] = res;
        offset += shape[i];
    }
```

```
shape:      [2, 3, 1]
mat_cols:   [0,   1,   1,   2,   3,   3 ]
mat_vals:   [2.0, -1.0, -1.0, 2.0, -1.0, 3.0]
vct:        [1.0,  2.0,  3.0, 4.0]
```

---

# Sparse Matrix-Vector Multiplication: Pseudocode

```
smvMul( mat_cols: [N]int, mat_vals: [N]real
        , vct: [vct_len], shape: [num_rows]int ) : [num_rows]real =
    real res_vct[num_rows];
    int offset = 0;
    for(int i=0; i < num_rows; i++) {
        real res = 0;
        for(int j=offset; j < offset+shape[i]; j++) {
            int vct_ind = mat_cols[j];
            int mat_elm = mat_vals[j];
            res += mat_elm * vct[vct_ind];
        }
        res_vct[i] = res;
        offset += shape[i];
    }
```

```
shape:      [2, 3, 1]
mat_cols:   [0,   1,   1,   2,   3,   3 ]
mat_vals:   [2.0, -1.0, -1.0, 2.0, -1.0, 3.0]
vct:        [1.0,  2.0,  3.0, 4.0]
```

-----

Result:

```
[ 2.0*1.0 + (-1.0)*2.0           = 0.0
, (-1.0)*2.0 + 2.0*3.0 + (-1.0)*4.0 = 0.0
, 3.0*4.0                        = 12.0
]
```

# Sparse Matrix-Vector Multiplication

Imperative Pseudocode:

```
smvMul( mat_cols: [N]int, mat_vals: [N]real
        , vct: [vct_len], shape: [num_rows]int ) : [num_rows]real =
    real res_vct[num_rows];
    int offset = 0;
    for(int i=0; i < num_rows; i++) {
        real res = 0;
        for(int j=offset; j < offset+shape[i]; j++) {
            int vct_ind = mat_cols[j];
            int mat_elm = mat_vals[j];
            res += mat_elm * vct[vct_ind];
        }
        res_vct[i] = res;
        offset += shape[i];
    }
```

We would like to:

- ▶ exploit all parallelism, **but the inner parallelism is irregular!**
- ▶ the inner loop:

# Sparse Matrix-Vector Multiplication

Imperative Pseudocode:

```
smvMul( mat_cols: [N]int, mat_vals: [N]real
        , vct: [vct_len], shape: [num_rows]int ) : [num_rows]real =
    real res_vct[num_rows];
    int offset = 0;
    for(int i=0; i < num_rows; i++) {
        real res = 0;
        for(int j=offset; j < offset+shape[i]; j++) {
            int vct_ind = mat_cols[j];
            int mat_elm = mat_vals[j];
            res += mat_elm * vct[vct_ind];
        }
        res_vct[i] = res;
        offset += shape[i];
    }
```

We would like to:

- ▶ exploit all parallelism, **but the inner parallelism is irregular!**
- ▶ the inner loop: looks like a map-reduce composition, but its size is variant to the outer map;
- ▶ Hint: compute flags + use segmented scan + extract last elements!

# Computing Flag Array from Shape Array

This works even in the presence of empty rows!

```
mkFlagArray (aoa_shp: [m]i32) : []i32 =    aoa_shp=[0,3,1,0,4,2,0]
  shp_rot = map (\i->if i==0 then 0
                  else aoa_shp[i-1]
                  , [0,...,m-1]);
  shp_scn = scan( (+), 0, shp_rot );
  aoa_len = shp_scn[m-1] + aoa_shp[m-1];
  shp_ind = map2 (\shp ind ->
                  if shp==0 then -1
                  else ind
                  , aoa_shp, shp_scn);
  return
    scatter( replicate(aoa_len, 0)
            , shp_ind
            , replicate(m, 1) );
```

# Computing Flag Array from Shape Array

This works even in the presence of empty rows!

```
mkFlagArray (aoa_shp: [m]i32) : []i32 =  
  shp_rot = map (\i->if i==0 then 0  
                  else aoa_shp[i-1]  
                  , [0,...,m-1]);  
  shp_scn = scan( (+), 0, shp_rot );  
  aoa_len = shp_scn[m-1] + aoa_shp[m-1];  
  shp_ind = map2 (\shp ind ->  
                  if shp==0 then -1  
                  else ind  
                  , aoa_shp, shp_scn);  
  return  
    scatter( replicate(aoa_len, 0)  
            , shp_ind  
            , replicate(m, 1) );
```

aoa\_shp=[0,3,1,0,4,2,0]  
shp\_rot=[0,0,3,1,0,4,2]

# Computing Flag Array from Shape Array

This works even in the presence of empty rows!

```
mkFlagArray (aoa_shp: [m]i32) : []i32 =  
  shp_rot = map (\i->if i==0 then 0  
                  else aoa_shp[i-1]  
                  , [0,...,m-1]);  
  shp_scn = scan( (+), 0, shp_rot );  
  aoa_len = shp_scn[m-1] + aoa_shp[m-1];  
  shp_ind = map2 (\shp ind ->  
                  if shp==0 then -1  
                  else ind  
                  , aoa_shp, shp_scn);  
  return  
    scatter( replicate(aoa_len, 0)  
            , shp_ind  
            , replicate(m, 1) );
```

aoa\_shp=[0,3,1,0,4,2,0]  
shp\_rot=[0,0,3,1,0,4,2]  
shp\_scn=[0,0,3,4,4,8,10]



# Computing Flag Array from Shape Array

This works even in the presence of empty rows!

```
mkFlagArray (aoa_shp: [m]i32) : []i32 =  
  shp_rot = map (\i->if i==0 then 0  
                  else aoa_shp[i-1]  
                  , [0,...,m-1]);  
  shp_scn = scan( (+), 0, shp_rot );  
  aoa_len = shp_scn[m-1] + aoa_shp[m-1];  
  shp_ind = map2 (\shp ind ->  
                  if shp==0 then -1  
                  else ind  
                  , aoa_shp, shp_scn);  
  return  
    scatter( replicate(aoa_len, 0)  
            , shp_ind  
            , replicate(m, 1) );
```

aoa\_shp=[0,3,1,0,4,2,0]  
shp\_rot=[0,0,3,1,0,4,2]  
shp\_scn=[0,0,3,4,4,8,10]  
aoa\_len= 10  
shp\_ind=

# Computing Flag Array from Shape Array

This works even in the presence of empty rows!

```
mkFlagArray (aoa_shp: [m]i32) : []i32 =  
  shp_rot = map (\i->if i==0 then 0  
                  else aoa_shp[i-1]  
                  , [0,...,m-1]);  
  shp_scn = scan( (+), 0, shp_rot );  
  aoa_len = shp_scn[m-1] + aoa_shp[m-1];  
  shp_ind = map2 (\shp ind ->  
                  if shp==0 then -1  
                  else ind  
                  , aoa_shp, shp_scn);  
  return  
    scatter( replicate(aoa_len, 0)  
            , shp_ind  
            , replicate(m, 1) );
```

```
aoa_shp=[0,3,1,0,4,2,0]  
shp_rot=[0,0,3,1,0,4,2]  
  
shp_scn=[0,0,3,4,4,8,10]  
aoa_len= 10  
  
shp_ind=  
  [-1,0,3,-1,4,8,-1]  
scatter
```

# Computing Flag Array from Shape Array

This works even in the presence of empty rows!

```
mkFlagArray (aoa_shp: [m]i32) : [m]i32 =  
  shp_rot = map (\i->if i==0 then 0  
                  else aoa_shp[i-1]  
                  , [0,...,m-1]);  
  shp_scn = scan( (+), 0, shp_rot );  
  aoa_len = shp_scn[m-1] + aoa_shp[m-1];  
  shp_ind = map2 (\shp ind ->  
                  if shp==0 then -1  
                  else ind  
                  , aoa_shp, shp_scn);  
  return  
    scatter( replicate(aoa_len, 0)  
            , shp_ind  
            , replicate(m, 1) );
```

```
aoa_shp=[0,3,1,0,4,2,0]  
shp_rot=[0,0,3,1,0,4,2]  
  
shp_scn=[0,0,3,4,4,8,10]  
aoa_len= 10  
  
shp_ind=  
  [-1,0,3,-1,4,8,-1]  
scatter  
  [0,0,0,0,0,0,0,0,0,0]  
  [-1,0,3,-1,4,8,-1]  
  [1,1,1,1,1,1,1]  
  -----  
  [1,0,0,1,1,0,0,0,1,0]
```

# Flat-Parallel Sparse Matrix-Vector Mult

Idea: use segmented scan, then extract last element in every row!

We assume no rows of length 0, and also inline flag computation.

```
smvMul( mat_cols: [N]int           // [0, 1, 1, 2, 3, 3]
        , mat_vals: [N]real        // [2.0, -1.0, -1.0, 2.0, -1.0, 3.0]
        , vct: [vct_len]           // [1.0, 2.0, 3.0, 4.0]
        , shape: [num_rows]int )   // [2, 3, 1]
    : [num_rows]real =              // Expected: [0.0, 0.0, 12.0]
    shpscn= scan( (+), 0, shape );   // [2, 5, 6]
    len    = shpscn[num_rows-1];
    ones   = replicate(num_rows, 1); // [1, 1, 1]
    indm1  = map(λ i -> if i==0 then 0
                  else shpscn[i-1]
                  , [0...num_rows]); // [0, 2, 5]
    flags  = scatter( replicate(len, 0) // [0, 0, 0, 0, 0, 0]
                     , indm1
                     , ones );         // [1, 0, 1, 0, 0, 1]
    prods  = map2(λ col v ->          // [2.0, -2.0, -2.0, 6.0, -4.0, 12.0]
                  v * vct[col]
                  , mat_cols, mat_vals );
    scnarr= sgmScaninc( (+), 0
                       , flags, prods); // [2.0, 0.0, -2.0, 4.0, 0.0, 12.0]
```

# Flat-Parallel Sparse Matrix-Vector Mult

Idea: use segmented scan, then extract last element in every row!

We assume no rows of length 0, and also inline flag computation.

```
smvMul( mat_cols: [N]int           // [0, 1, 1, 2, 3, 3]
        , mat_vals: [N]real        // [2.0, -1.0, -1.0, 2.0, -1.0, 3.0]
        , vct: [vct_len]           // [1.0, 2.0, 3.0, 4.0]
        , shape: [num_rows]int )   // [2, 3, 1]
    : [num_rows]real =              // Expected: [0.0, 0.0, 12.0]
    shpscn= scan( (+), 0, shape );   // [2, 5, 6]
    len    = shpscn[num_rows-1];
    ones   = replicate(num_rows, 1); // [1, 1, 1]
    indm1  = map(λ i -> if i==0 then 0
                  else shpscn[i-1]
                  , [0..num_rows]); // [0, 2, 5]
    flags  = scatter( replicate(len, 0) // [0, 0, 0, 0, 0, 0]
                     , indm1
                     , ones );         // [1, 0, 1, 0, 0, 1]
    prods  = map2(λ col v ->          // [2.0, -2.0, -2.0, 6.0, -4.0, 12.0]
                  v * vct[col]
                  , mat_cols, mat_vals );
    scnarr = sgmScaninc( (+), 0
                        , flags, prods); // [2.0, 0.0, -2.0, 4.0, 0.0, 12.0]
    return                                     // [scnarr[1], scn_arr[4], scn_arr[5]]
        map(λ ind -> scnarr[ind-1];
            , shpscn );                       // [0.0, 0.0, 12.0]
```

## Exercise 4: Sparse Matrix-Vector Mult in OpenCL

- ▶ the CPU stub is in  
Day3-Exercises/ScanApps/spMatVecMult.h,  
function runSpMatVectMul; nothing to do there! You only  
need to implement several kernels in scanapps.cl.
- (1) `shpscn= scan( (+), 0, shape );`  
is already implemented!
- (2) Implement kernel `mkFlagsSpMVM`. Array `flags` has been  
already zeroed; `shape_scn` is `shpscn` in the previous slide.
- (3) Implement kernel `mulPhaseSpMVM`.  
`mat_ind` is `mat_cols` from the previous slide.
- (4) `scnarr= sgmScaninc( (+), 0, flags, prods);`  
is already implemented!
- (5) Implement kernel `getLastSpMVM`.  
`shape_scn` and `sgm_mat` are `shpscn` and `scnarr` in  
previous slide, respectively, and `out` is the result.

## Course Contents

### Map and Reduce

- Types, Semantics and Properties

- Efficient Sequentialization: Work, Depth, and Brent Lemma

- Maximum Segment Sum Problem (MSSP)

- Exercise 1: OpenCL Implementation of MSSP

### Scan and Applications

- Scan: Type, Semantics, Asymptotic, Parallel Implementation

- Exercise 2: Intra-Wave Inclusive Scan Implementation

- Other Data-Parallel Operators: Scatter, Partition

- Exercise 3, Step 2: Partition2 OpenCL Implementation

### Segmented Scan and Applications

- Segmented Scan: Type, Semantics, Asymptotic, Impleme

- Application: Sparse Matrix-Vector Multiplication

- Exercise 4: Sparse Matrix-Vector Mult in OpenCL

## Optimized Partition2

## Back To Partition2

```
partition2(X: [n]i32): (i32, [n]i32) =  
  cs = map(even, X);  
  tfs= map (\c->if c then 1 else 0  
            , cs);  
  isT= scan( (+), 0, tfs );  
  i   = isT[n-1];  
  
  ffs= map (\c->if c then 0 else 1  
            , cs );  
  tmps = scan( (+), 0, ffs );  
  isF= map (\t -> t+i, tmps);  
  
  inds=map3(\(c,iT,iF) ->  
            if c then iT-1 else iF-1  
            , cs, isT, isF);  
  res = scatter(scratch n, inds, X);  
  return (i, res);
```

```
X    = [5, 4, 2, 10,3, 7, 8]  
n    = 7  
cs   = [F, T, T, T, F, F, T]  
tfs  = [0, 1, 1, 1, 0, 0, 1]  
  
isT  = [0, 1, 2, 3, 3, 3, 4]  
i    = 4  
  
ffs  = [1, 0, 0, 0, 1, 1, 0]  
tmps = [1, 1, 1, 1, 2, 3, 3]  
isF  = [5, 5, 5, 5, 6, 7, 7]  
  
inds = [4, 0, 1, 2, 5, 6, 3]  
X     = [5, 4, 2, 10,3, 7, 8]  
Res   = [_, _, _, _, _, _, _]  
-----  
Res   = [4, 2, 10, 8,5, 3, 7]
```

This shows you how to think about it, but a direct implementation is not very efficient in practice.



## Partition2: Taking Inspiration From Scan

```
bool even(int v) { return (bool)(1 -(v&1)); }  
partition2(even, [5, 4, 2, 10, 3, 7, 8] ) should result in  
          (4, [4, 2, 10, 8, 5, 3, 7])
```

A scan implementation is typically split into three kernels:

- (1) reduce at workgroup level (with sequential chunking);
- (2) scan the results of (1) in one workgroup; this is typically negligible, as it is applied on few elements;
- (3) perform the workgroup level scan while starting from the result of the previous workgroup.

## Partition2: Deriving an Optimized Implementation

```
bool even(int v) { return (bool)(1 -(v&1)); }  
partition2(even, [5, 4, 2, 10, 3, 7, 8] ) should result in  
          (4, [4, 2, 10, 8, 5, 3, 7])
```

- (I) map-reduce at workgroup level; operator is commutative!
  - (a)  $X = \text{map}(\lambda(e) \rightarrow \text{if } \text{pred}(e) \text{ then } (1,0) \text{ else } (0,1), A);$
  - (b)  $\text{reduce}^{wkgp}(\lambda((a1,b1),(a2,b2)) \rightarrow (a1+b2, b1+b2)), (0,0), X);$
- (II) scan the results of (I) in one workgroup; after this step we know the global offsets of the elements that succeed/fail at the start of each workgroup.
- (III) Do all the rest in one kernel:
  - (1) Copy ELEMS\_PER\_THREAD from global to local memory, coalesced (ToDo).
  - (2) Save them in per-thread register memory (chunk) +  
Reduce sequentially each thread's ELEMS\_PER\_THREAD elems in `tf`; (ToDo).
  - (3) Perform workgroup-level scan on the `tf` of each thread (DONE).
  - (4) The number of elements that succeed/fail per workgroup (`num_ts/num_fs`) is the last scanned element (DONE).
  - (5) Save in `tf` the result of the previous thread (DONE).
  - (6) Redo computation (2) but instead place the elements from register to local memory in the order in which they would occur if the partition was done at workgroup level; remember to add `num_ts` to the failed-elems index (ToDo).
  - (7) Finally copy-out from local to global memory in coalesced fashion (DONE).

## Last Exercise: Optimized Partition2

- ▶ CPU stub: in file `Day3-Exercises/Partition2Opt/partition2.h`, function `runPartition`; that is where the OpenCL kernels are called, but nothing to do there.
- ▶ Open in file `Day3-Exercises/Partition2Opt/kernels.cl`, kernel `scanPhaseKer`. Follow the instructions, which correspond to implementing the `ToDo` items from the previous slide.
- ▶ The performance of the optimized `partition2` should come close to that of inclusive scan!