

Parallel Algorithms Patterns: Scan

COMPUTAÇÃO DE ALTO DESEMPENHO 2018/2019

HERVÉ PAULINO

SLIDES ADAPTED FROM

“STRUCTURED PARALLEL PROGRAMMING - CIS 410/510, UNIVERSITY OF OREGON”

REZAUL A. CHOWDHURY - CSE 613: PARALLEL PROGRAMMING



Bibliography

Section 3.3.5 of Structured Parallel Programming, Michael McCool, Arch D. Robison and James Reinders. Morgan Kaufmann (2012)

Parallel Prefix Sum with CUDA, Mark Harris,
<https://www.mimuw.edu.pl/~ps209291/kgkp/slides/scan.pdf>

Scan or Prefix Sum

Scan: computes all partial reduction of a collection

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
5	3	7	1	3	6	2	4

\oplus = binary addition

5	8	15	16	19	25	27	31
s_1	s_2	s_3	s_4	s_5	s_6	s_7	s_8

For every output in a collection, a reduction of the input up to that point is computed

If the function being used is associative, the scan can be parallelized

Parallelizing a scan is not obvious at first, because of dependencies to previous iterations in the serial loop

A parallel scan will require more operations than a serial version

Exclusive Scan

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
5	3	7	1	3	6	2	4

\oplus = binary addition

0	5	8	15	16	19	25	27	31
---	---	---	----	----	----	----	----	----



Identity

Contraction Technique

1. **Reduce:** reduce the original problem to a smaller problem
2. **Conquer:** solve the smaller problem (often recursively)
3. **Expand:** use the solution to the smaller problem to obtain a solution for the original larger problem

Contraction

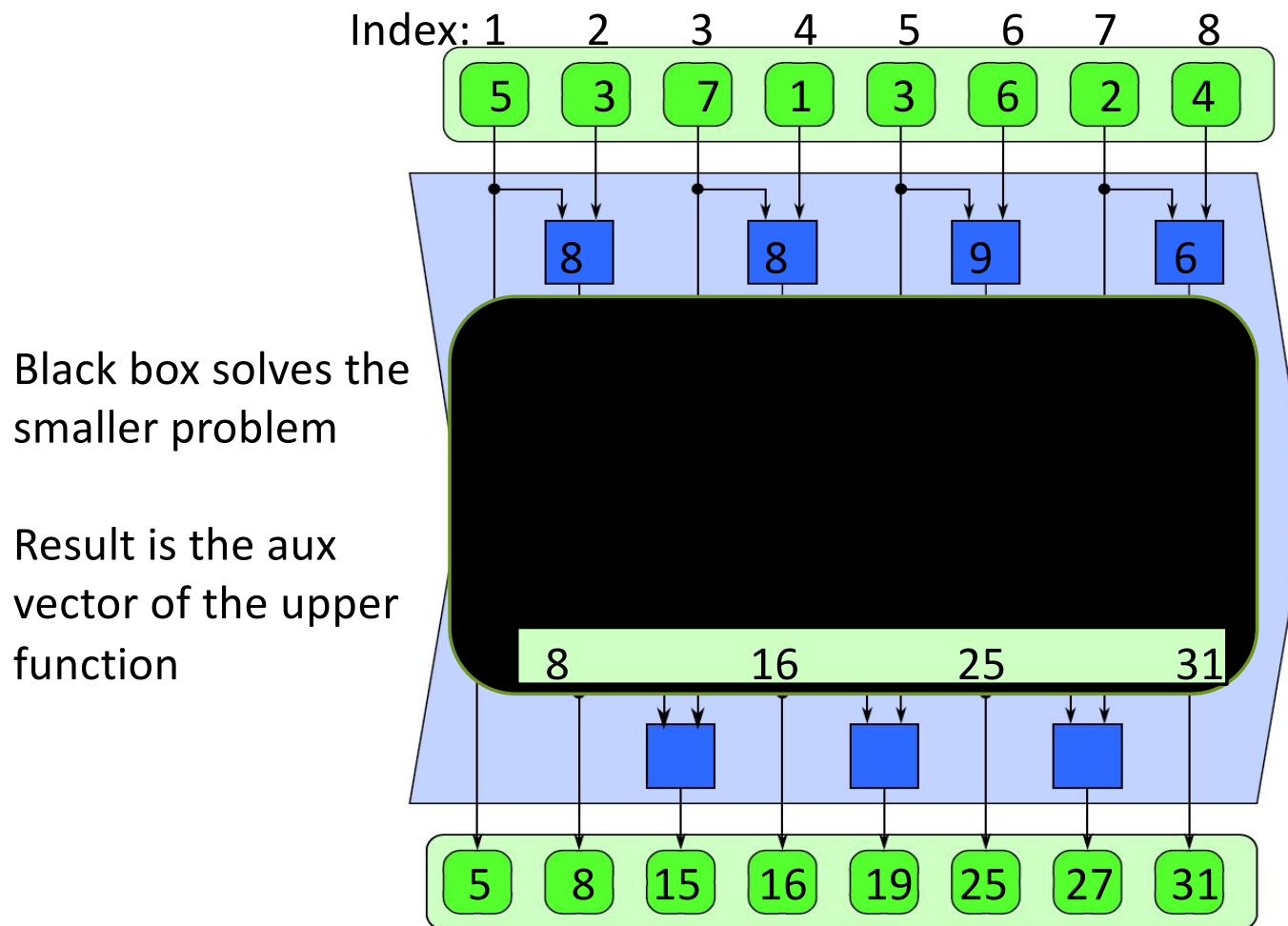
Scan or Prefix Sum in CPU

```
prefix_sum(array) {  
  if array.size() == 1 then result[1] = array[1]  
  else  
    parallel for i = 1 to array.size()/2 do  
      aux[i] = array[2*i-1] + array[2*i]  
    aux = prefix_sum(aux)  
    parallel for i = 1 to array.size() do  
      if i = 1 then result[1] = array[1]  
      else if i = even then result[i] = aux[i/2]  
      else result[i] = aux[(i-1)/2] + array[i]  
    return result  
}
```

Solve problem by processing
the result of solving a smaller
instance of the problem

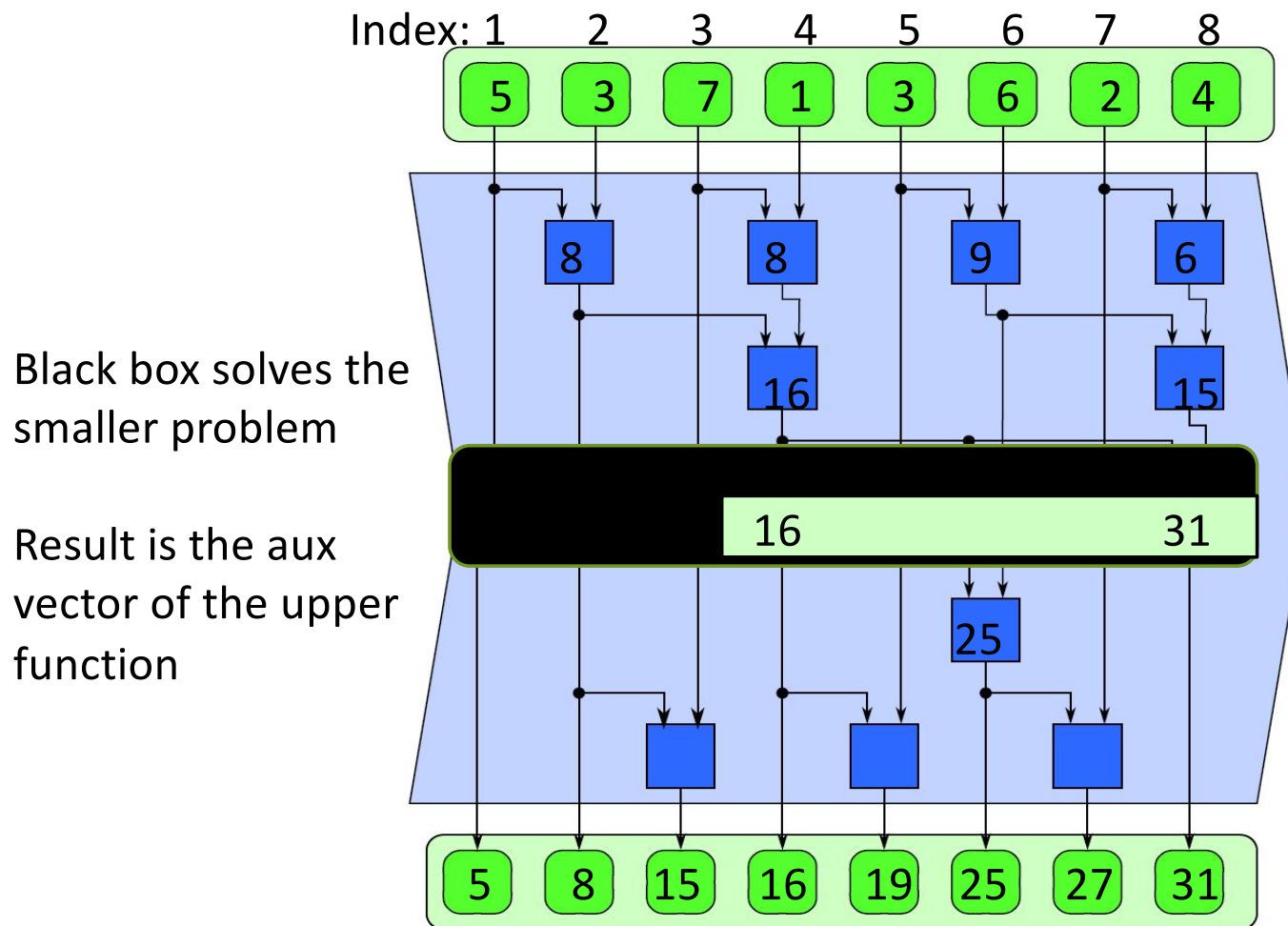
Contraction

Scan or Prefix Sum



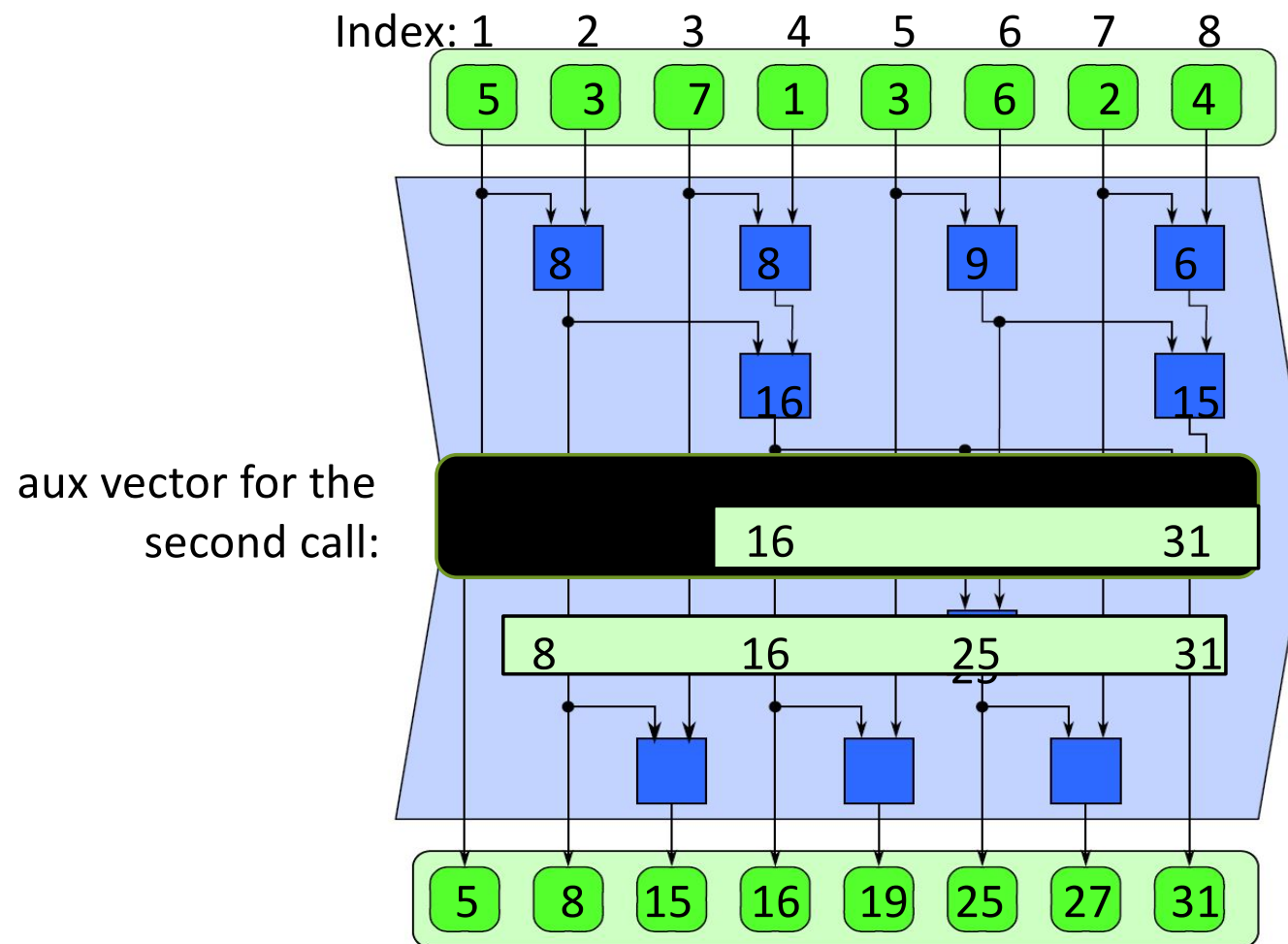
Contraction

Scan or Prefix Sum



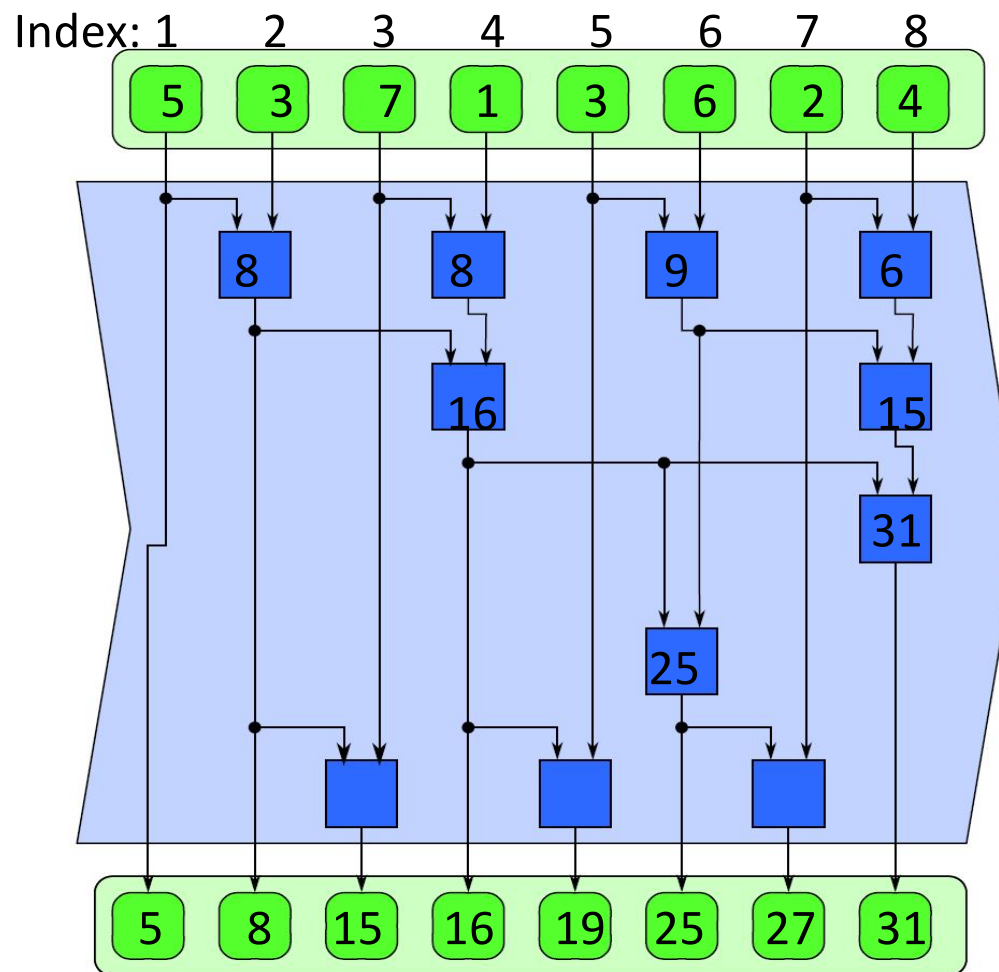
Contraction

Scan or Prefix Sum



Contraction

Scan or Prefix Sum



A GPU Implementation in a Single Block

Two phases:

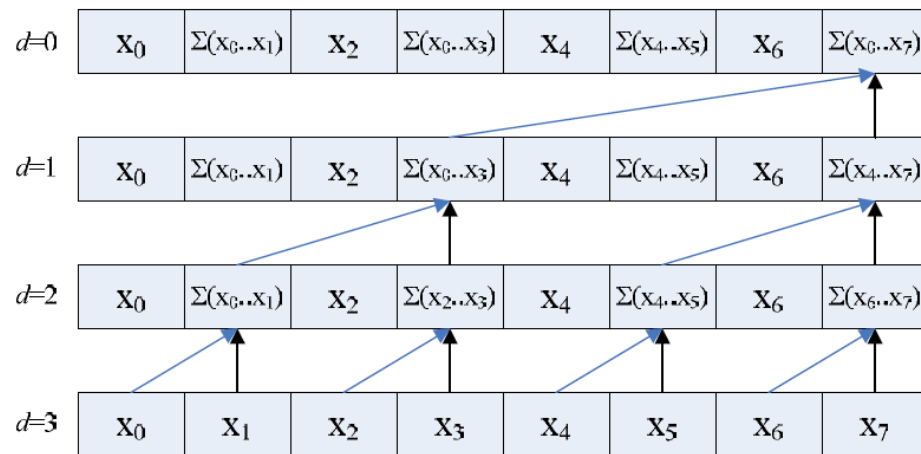
- up-sweep (reduce)
- down-sweep

Thread number: $\text{NUM_ELEMENTS} / 2$

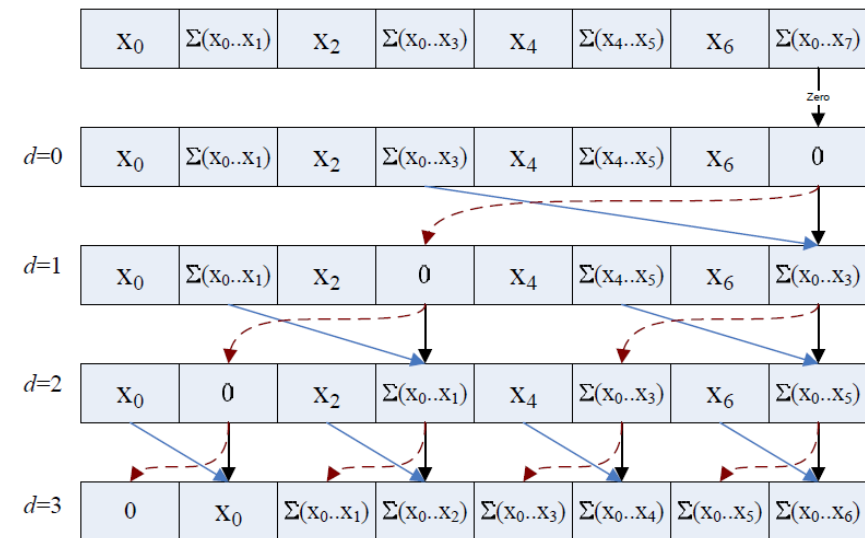
Shared memory size: NUM_ELEMENTS

A GPU Implementation in a Single Block

UP-SWEEP



DOWN SWEEP

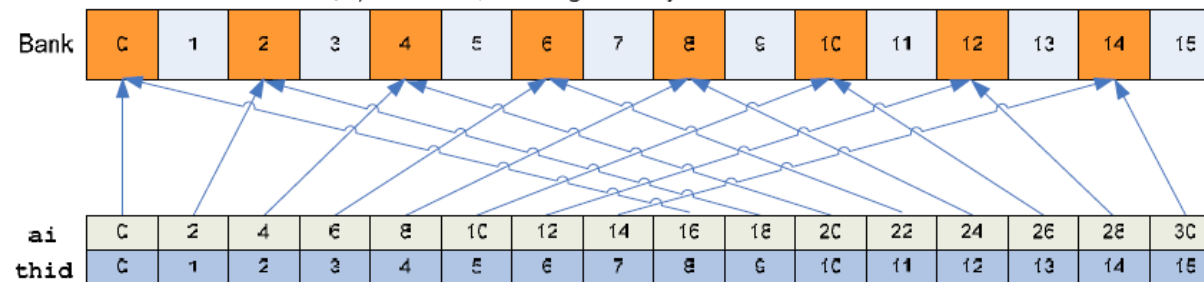


Avoiding Bank Conflicts

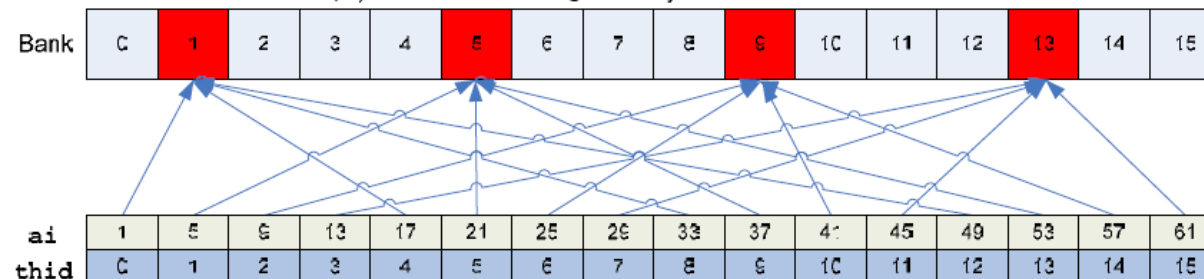
Addressing Without Padding

```
int ai = offset*(2*thid+1)-1;  
int bi = offset*(2*thid+2)-1;  
temp[bi] += temp[ai];
```

Offset = 1: Address (ai) stride is 2, resulting in 2-way bank conflicts



Offset = 2: Address (ai) stride is 4, resulting in 4-way bank conflicts

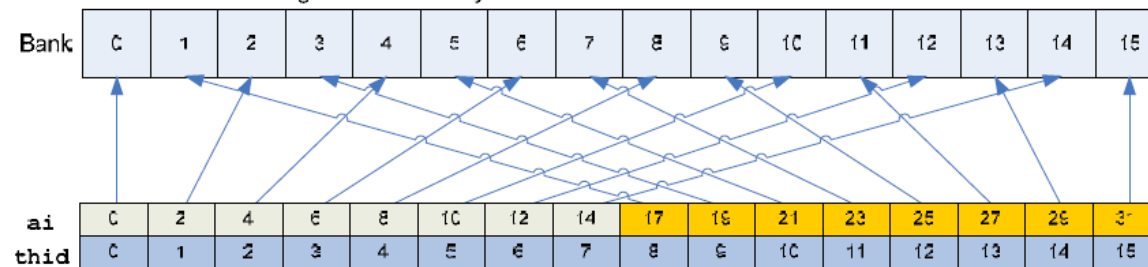


Avoiding Bank Conflicts

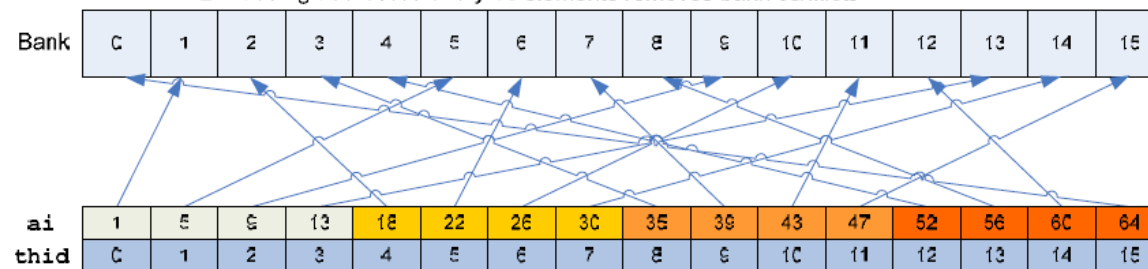
Addressing With Padding

```
int ai = offset*(2*thid+1)-1;
int bi = offset*(2*thid+2)-1;
ai += ai / NUM_BANKS;
bi += bi / NUM_BANKS;
temp[bi] += temp[ai];
```

offset = 1: Padding addresses every 16 elements removes bank conflicts



offset = 2: Padding addresses every 16 elements removes bank conflicts



Padding increment:

0	1	2	3
---	---	---	---

Arrays of Arbitrary Size

