# Stencil Pattern

Parallel Computing

CIS 410/510

Department of Computer and Information Science

**O** | UNIVERSITY OF OREGON

# *Table of Contents*

❑ What is the stencil pattern?

❑ Implementing stencil with shift

❑ Stencil and cache optimizations
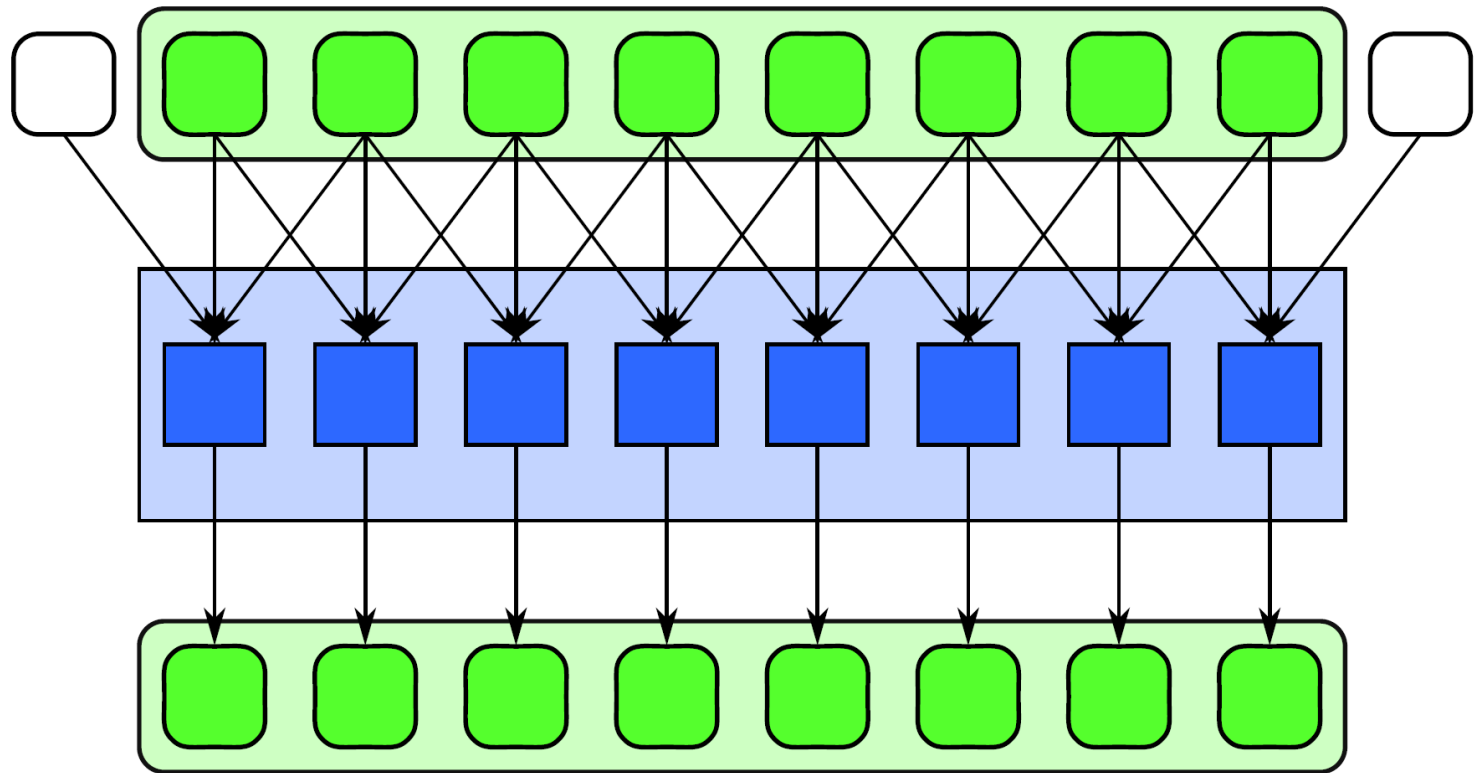
❑ Stencil and communication optimizations

❑ Recurrence

# *Table of Contents*

❑ What is the stencil pattern?

❑ Implementing stencil with shift

❑ Stencil and cache optimizations
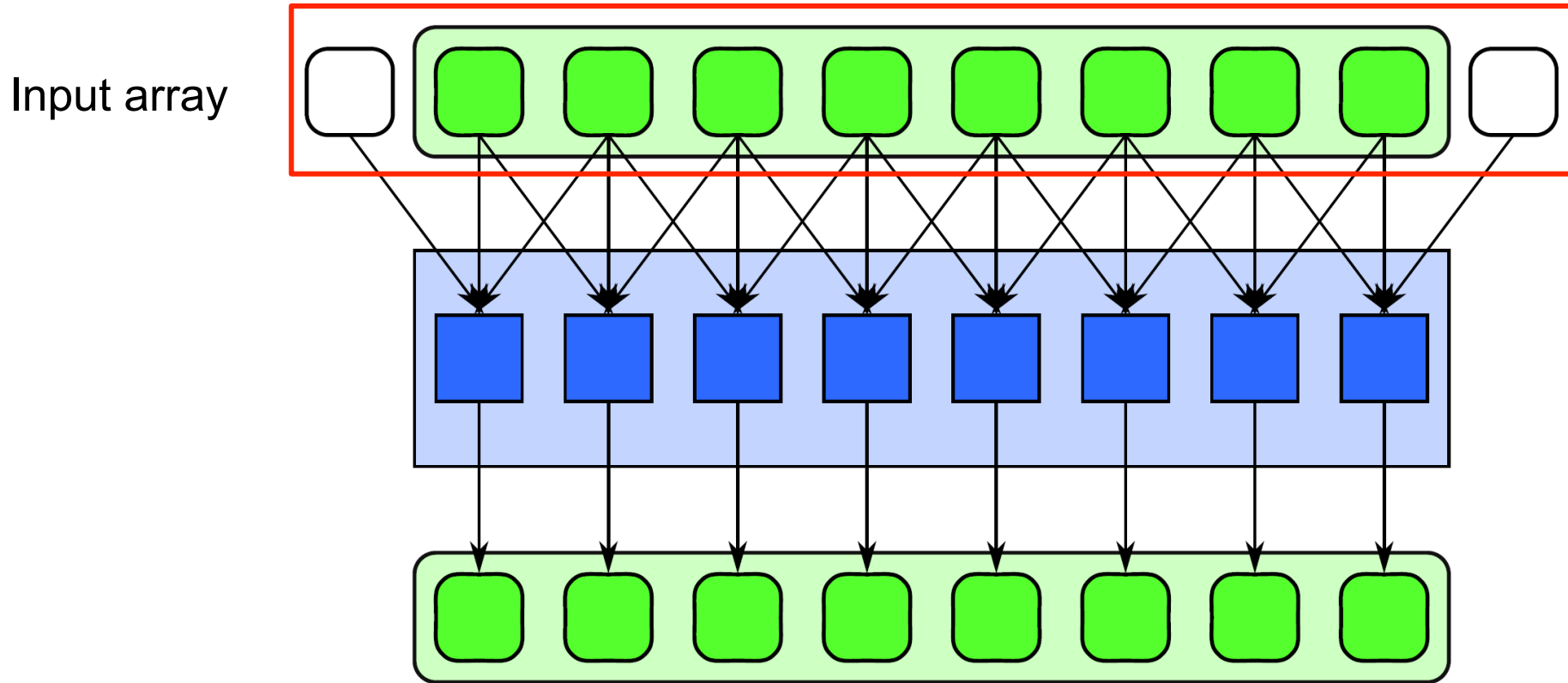
❑ Stencil and communication optimizations

❑ Recurrence

# *What is the stencil pattern?*

❑ **Stencil**: A map where each output depends on a "neighborhood" of inputs. These inputs are a set of fixed offsets relative to the output position.

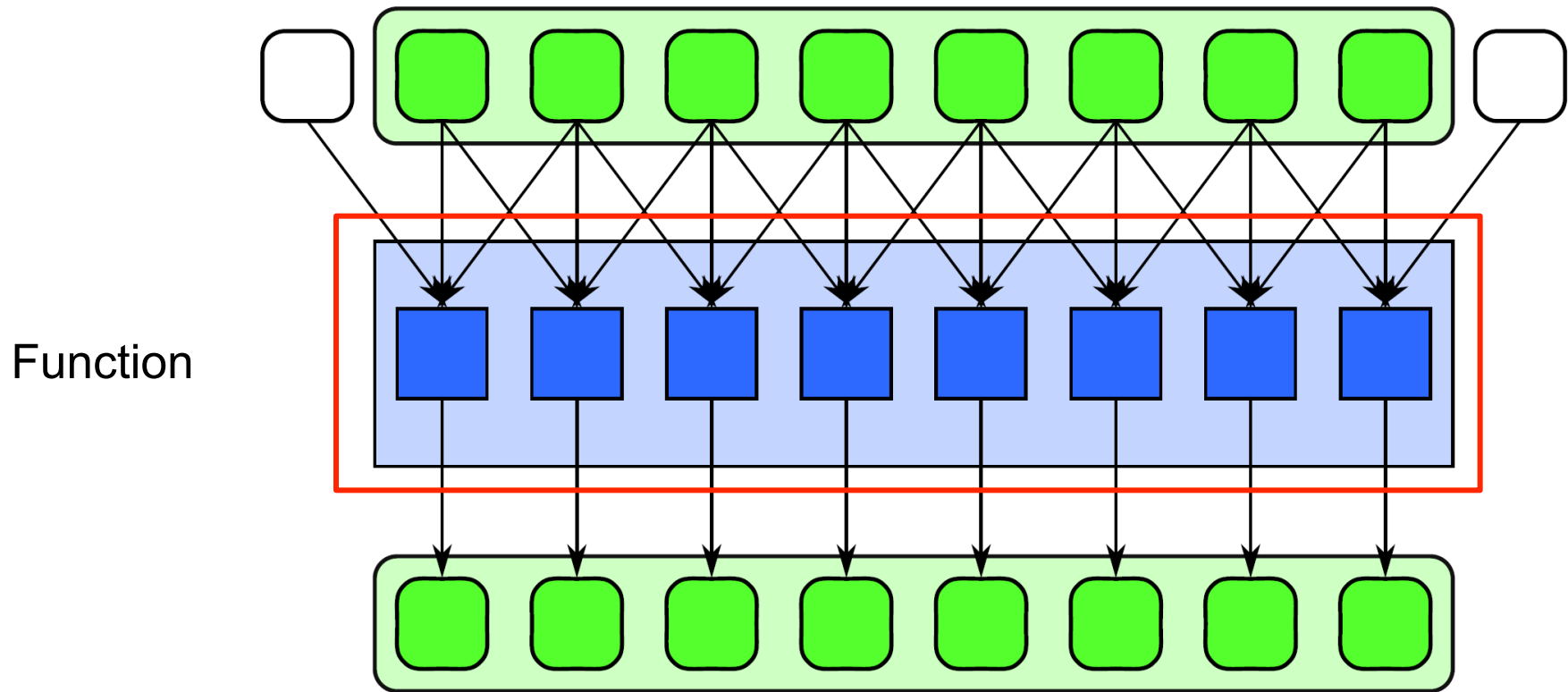❑ A stencil output is a function of a "neighborhood" of elements in an input collection
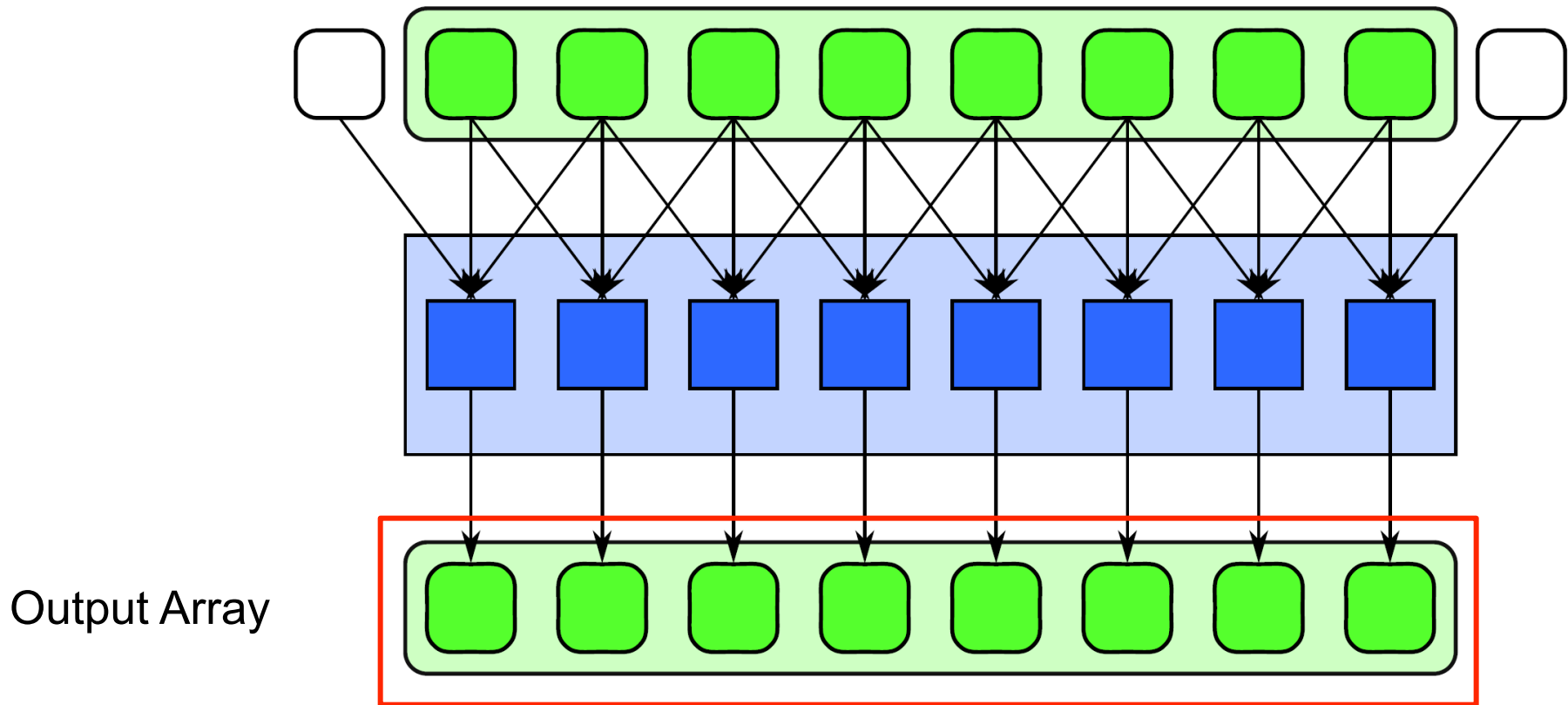
# *What is the stencil pattern?*

# *What is the stencil pattern?*

Input array

# *What is the stencil pattern?*

Function

# *What is the stencil pattern?*



Output Array

# *What is the stencil pattern?*

i-1    i    i+1

This stencil takes the i-1, i, i+1 elements…

# *What is the stencil pattern?*

Applies some function to them…

# *What is the stencil pattern?*

And outputs to the $i^{th}$ position of the output array

# *What is the stencil pattern?*

❑ Stencils can operate on one dimensional and multidimensional data

❑ Stencil neighborhoods can range from compact to sparse, square to cube, and anything else!

# *What is the stencil pattern?*

# Practice!

# What is the stencil pattern?

- Here is our array, A

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 9 | 7 | 0 |
| 0 | 6 | 4 | 0 |
| 0 | 0 | 0 | 0 |

# What is the stencil pattern?

- Here is our array, A

- Apply a stencil operation to the inner square of the form:

$(i,j) = \mathsf{avg}\big((i,j), (i-1,j), (i+1,j), (i,j-1), (i,j+1)\big)$

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 9 | 7 | 0 |
| 0 | 6 | 4 | 0 |
| 0 | 0 | 0 | 0 |

# What is the stencil pattern?

1) Average all blue squares

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 9 | 7 | 0 |
| 0 | 6 | 4 | 0 |
| 0 | 0 | 0 | 0 |

# What is the stencil pattern?

1)  Average all blue squares
2)  Store result in red square

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 9 | 7 | 0 |
| 0 | 6 | 4 | 0 |
| 0 | 0 | 0 | 0 |

UNIVERSITY OF OREGON

# What is the stencil pattern?

1) Average all blue squares
2) Store result in red square
3) Repeat 1 and 2 for all green squares

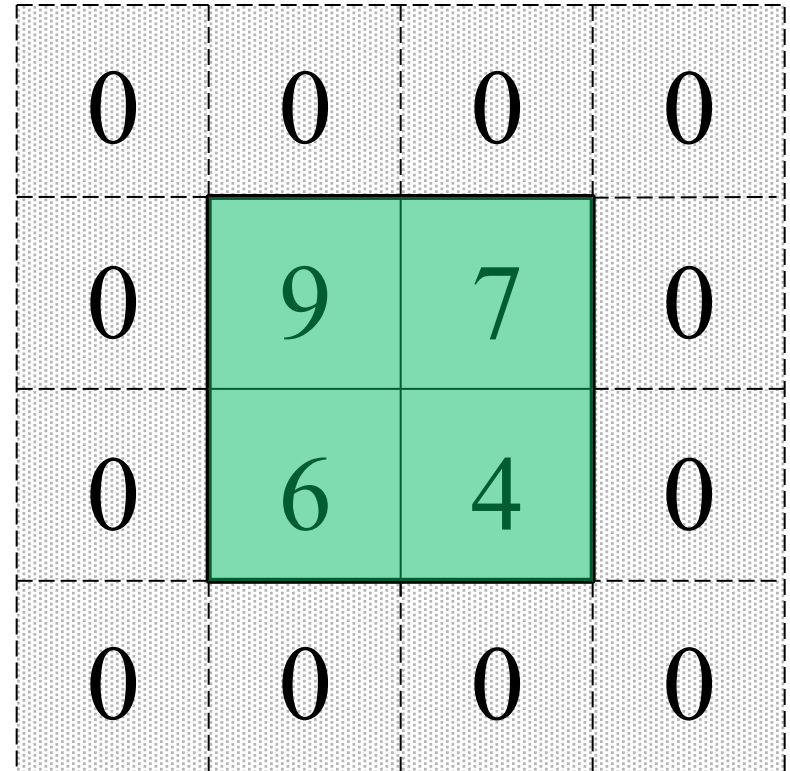| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 9 | 7 | 0 |
| 0 | 6 | 4 | 0 |
| 0 | 0 | 0 | 0 |

# What is the stencil pattern?

1) Average all blue squares
2) Store result in red square
3) Repeat 1 and 2 for all green squares

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 9 | 7 | 0 |
| 0 | 6 | 4 | 0 |
| 0 | 0 | 0 | 0 |

UNIVERSITY OF OREGON

# What is the stencil pattern?

- Output!

- How would we parallelize this?

| 4.4 | 4 |
|-----|-----|
| 3.8 | 3.4 |

# Serial Stencil Example (part 1)

```
1   template<
2       int NumOff,      // number of offsets
3       typename In,     // type of input locations
4       typename Out,    // type of output locations
5       typename F       // type of function/functor
6   >
7   void stencil(
8       int n,           // number of elements in data collection
9       const In a[],    // input data collection (n elements)
10      Out r[],         // output data collection (n elements)
11      In b,            // boundary value
12      F func,          // function/functor from neighborhood inputs to output
13      const int offsets[] // offsets (NumOffsets elements)
14  ) {
```

# Serial Stencil Example (part 2)

```
15    // array to hold neighbors
16    In neighborhood[NumOff];
17    // loop over all output locations
18    for (int i = 0; i < n; ++i) {
19        // loop over all offsets and gather neighborhood
20        for (int j = 0; j < NumOff; ++j) {
21            // get index of jth input location
22            int k = i+offsets[j];
23            if (0 <= k && k < n) {
24                // read input location
25                neighborhood[j] = a[k];
26            } else {
27                // handle boundary case
28                neighborhood[j] = b;
29            }
30        }
31        // compute output value from input neighborhood
32        r[i] = func(neighborhood);
33    }
34  }
```

UNIVERSITY
OF OREGON

# *Parallel Stencil Example*

# *Table of Contents*

❑ What is the stencil pattern?

❑ Implementing stencil with shift

❑ Stencil and cache optimizations
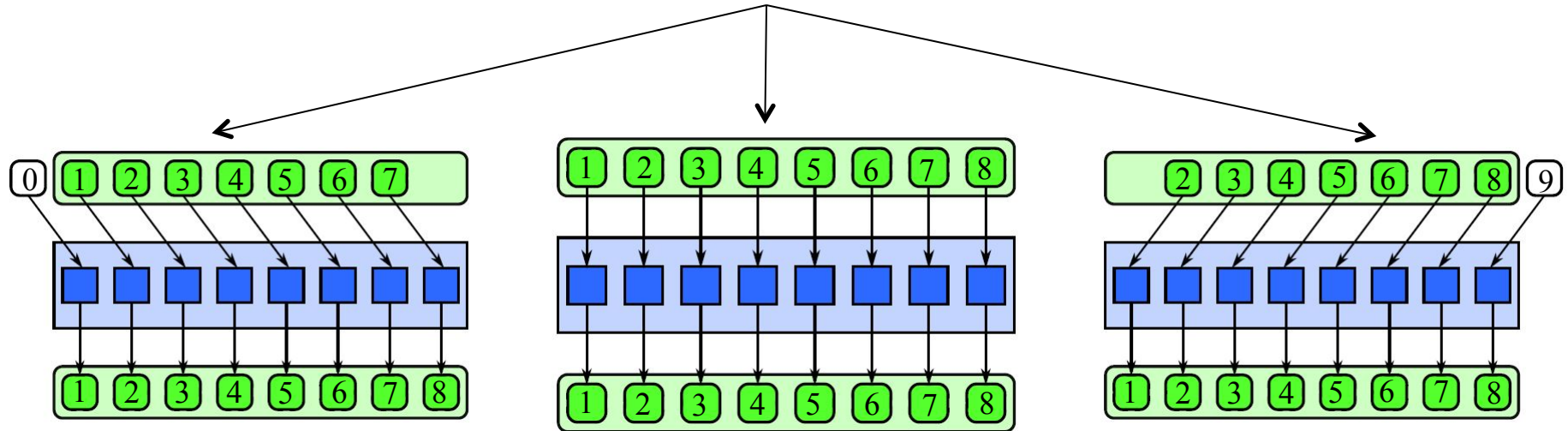
❑ Stencil and communication optimizations

❑ Recurrence

# *Implementing Stencil with Shift*

❑ One possible implementation of the stencil pattern includes shifting the input data

❑ For each offset in the stencil, we gather a new input vector by **shifting** the original input by the offset amount

# Implementing Stencil with Shift



All input arrays are derived from the same original input array

# *Implementing Stencil with Shift*

❑ This implementation is only beneficial for one dimensional stencils or the memory-contiguous dimension of a multidimensional stencil

❑ Memory traffic to external memory is not reduced with shifts

❑ But, shifts allow vectorization of the data reads, which may reduce the total number of instructions

# *Table of Contents*

❑ What is the stencil pattern?

❑ Implementing stencil with shift

❑ Stencil and cache optimizations

❑ Stencil and communication optimizations

❑ Recurrence

# *Stencil and Cache Optimizations*

❑ Assuming 2D array where rows are contiguous in memory…

  o Horizontally related data will tend to belong to the same cache line

  o Vertical offset accesses will most likely result in cache misses

# *Stencil and Cache Optimizations*

❑ Assigning rows to cores:

    o Maximizes horizontal data locality

    o Assuming vertical offsets in stencil, this will create redundant reads of adjacent rows from each core

❑ Assigning columns to cores:

    o Redundantly read data from same cache line

    o Create false sharing as cores write to same cache line

# *Stencil and Cache Optimizations*
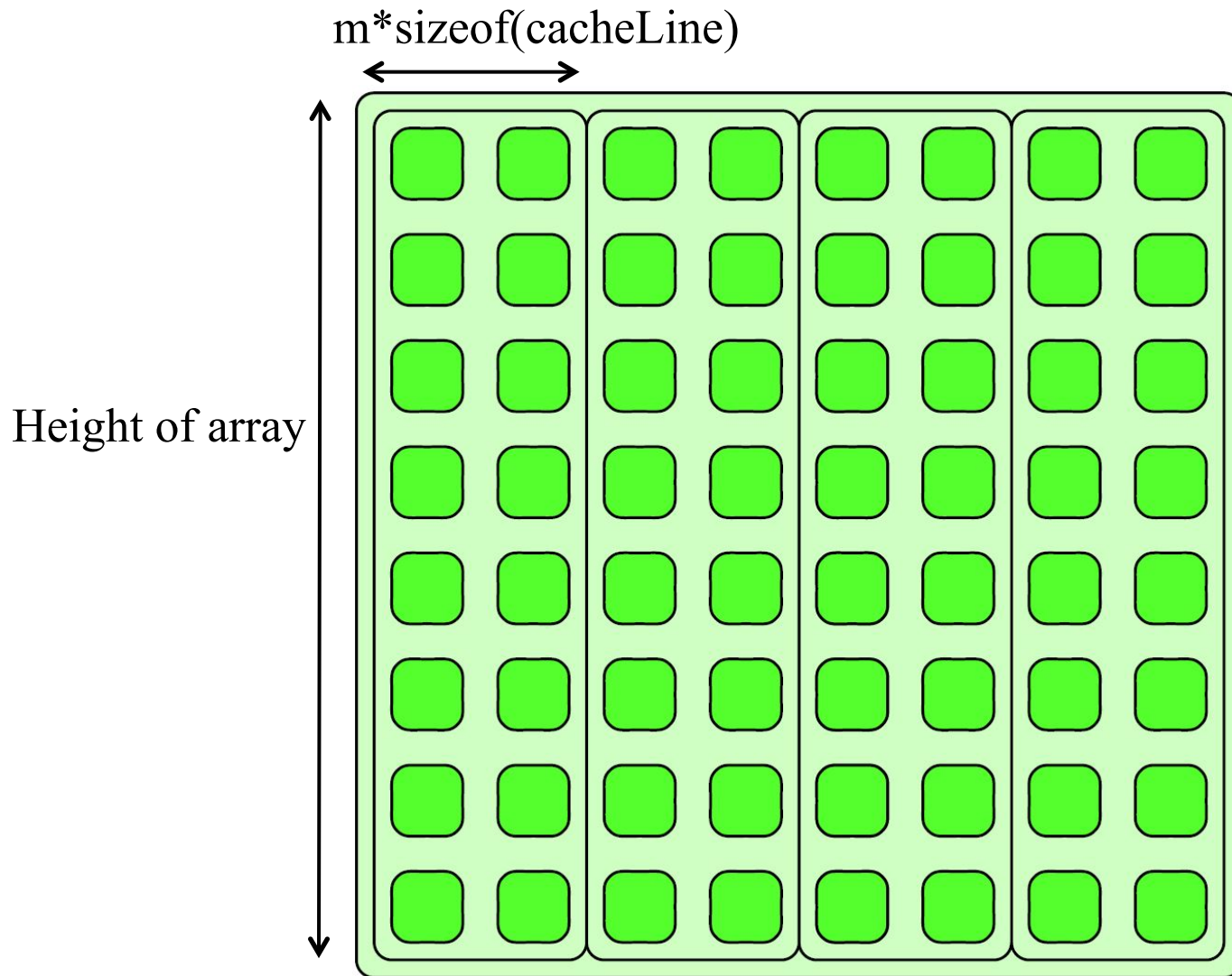
- ❑ Assigning "strips" to each core can be a better solution

- ❑ **Strip-mining**: an optimization in a stencil computation that groups elements in a way that avoids redundant memory accesses and aligns memory accesses with cache lines

# *Stencil and Cache Optimizations*

❑ A strip's size is a multiple of a cache line in width, and the height of the 2D array

❑ Strip widths are in increments of the cache line size so as to avoid false sharing and redundant reads

❑ Each strip is processed serially from top to bottom within each core

# *Stencil and Cache Optimizations*

m*sizeof(cacheLine)

Height of array

# *Table of Contents*

❑ What is the stencil pattern?

❑ Implementing stencil with shift

❑ Stencil and cache optimizations

❑ Stencil and communication optimizations

❑ Recurrence

# *Stencil and Communication Optimizations*

❏ When data is distributed, overlapping regions (or "**ghost cells**") must be explicitly communicated between nodes between loop iterations

❒ Darker cells are PE 0's ghost cells

❒ After first iteration of stencil computation

  ○ PE 1 & PE 2 must send their stencil results to PE 0

  ○ PE 0 can perform another iteration of stencil

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | PE 0<br>9 | PE 1<br>7 | 0 |
| 0 | PE 2<br>6 | PE 3<br>4 | 0 |
| 0 | 0 | 0 | 0 |

# *Stencil and Communication Optimizations*

❑ Generally better to replicate ghost cells in each local memory and swap after each iteration than to share memory

   o Fine-grained sharing can lead to increased communication cost

# *Stencil and Communication Optimizations*

❑ **Halo**: set of all ghost cells

❑ Halo must contain all neighbors needed for one iteration

❑ Larger halo (**deep halo**)
   o Less communications
   o More redundant local computation

❑ **Latency Hiding**: Compute interior of stencil while waiting for ghost cell updates

# *Table of Contents*

❑ What is the stencil pattern?

❑ Implementing stencil with shift

❑ Stencil and cache optimizations

❑ Stencil and communication optimizations

❑ Recurrence

# *Recurrence*

❑ What if we have several nested loops with data dependencies between them when doing a stencil computation?

# Recurrence

```
1   void my_recurrence(
2       size_t v,          // number of elements   vertically
3       size_t h,          // number of elements   horizontally
4       const float a[v][h], // input 2D array
5       float b[v][h]      // output 2D array (boundaries already   initialized )
6   ) {
7       for (int i=1; i<v; ++i)
8           for (int j=1; j<h; ++j)
9               b[i][j] = f(b[i-1][j], b[i][j-1], a[i][j]);
10  }
```

# *Recurrence*

```
1   void my_recurrence(
2       size_t v,           // number of elements   vertically
3       size_t h,           // number of elements   horizontally
4       const float a[v][h], // input 2D array
5       float b[v][h]       // output 2D array (boundaries already   initialized )
6   ) {
7       for (int i=1; i<v; ++i)
8           for (int j=1; j<h; ++j)
9               b[i][j] = f(b[i−1][j], b[i][j−1], a[i][j]);
10  }
```
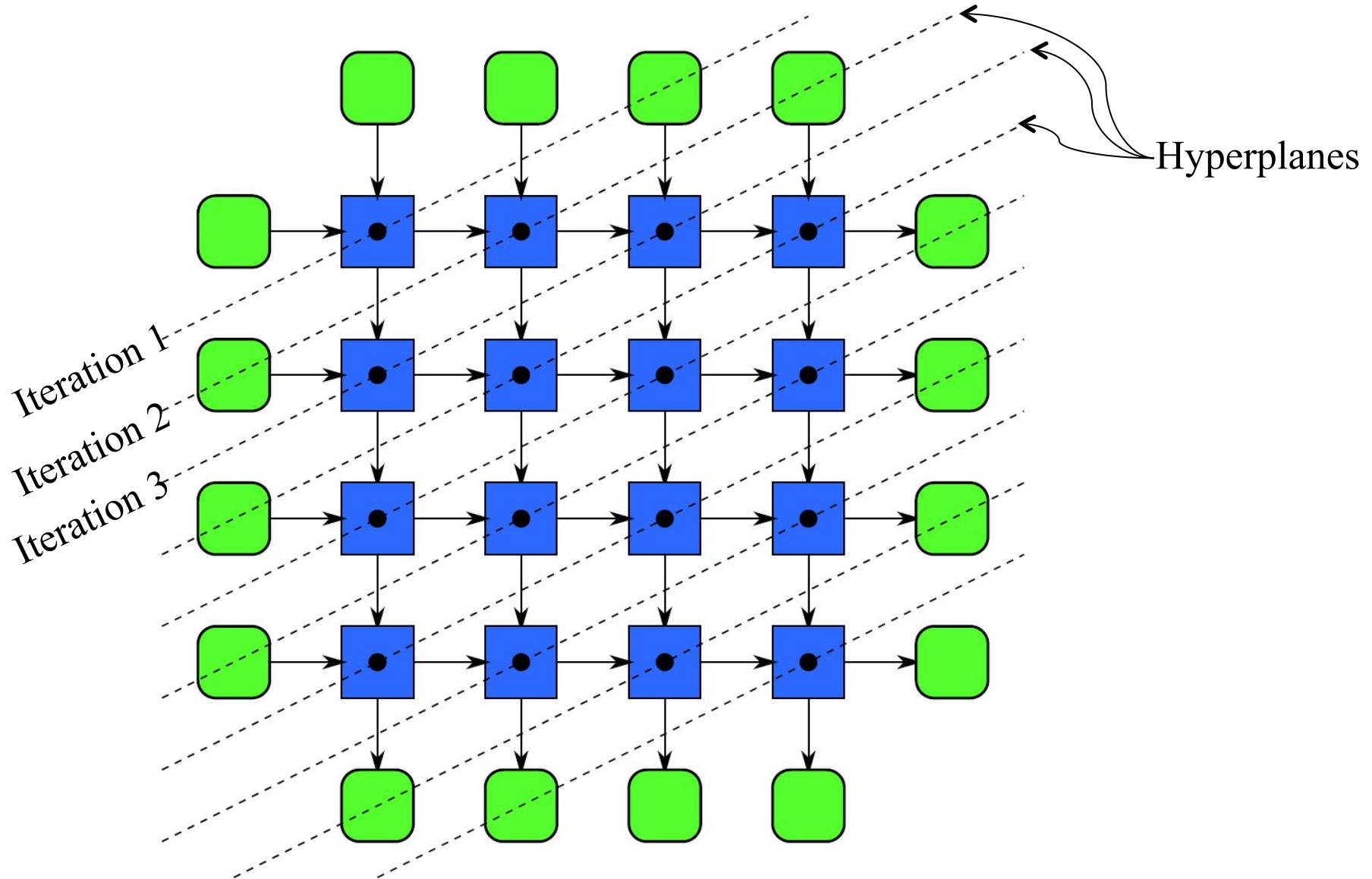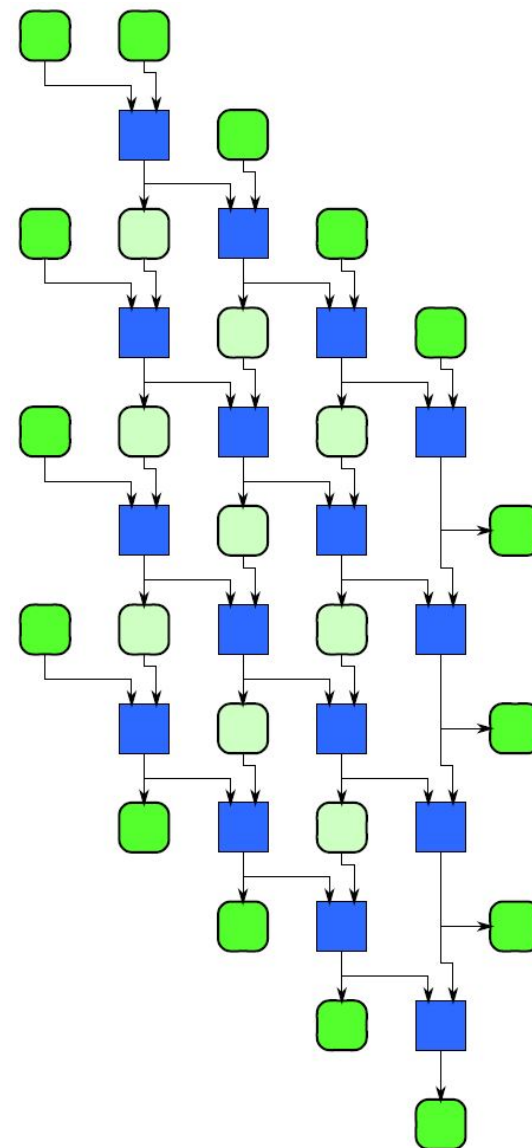
Data dependencies between loops

# *Recurrence*

❑ This can still be parallelized!

❑ Trick: find a plane that cuts through grid of intermediate results

    ○ Previously computed values on one side of plane

    ○ Values to still be computed on other side of plane

    ○ Computation proceeds perpendicular to plane through time (this is known as a sweep)

❑ This plane is called a *separating hyperplane*

# *Recurrence*



Hyperplanes

Iteration 1
Iteration 2
Iteration 3

# Recurrence

- Same grid of intermediate results

- Each level corresponds to a loop iteration

- Computation proceeds downward

# Example Implementation

# *Conclusion*

- ❑ Examined the stencil and recurrence pattern
    - o Both have a regular pattern of communication and data access
- ❑ In both patterns we can convert a set of offset memory accesses to shifts
- ❑ Stencils can use strip-mining to optimize cache use
- ❑ Ghost cells should be considered when stencil data is distributed across different memory spaces