Parallel Computing (CMSC416 / CMSC616)



Designing Parallel Programs

Abhinav Bhatele, Department of Computer Science



Reminders / Annoucements

- If you do not have a zaratan account, email: cmsc416@cs.umd.edu
- When emailing, please mention your course and section number:
 - Example: 416 / Section 0101
- Accomodations: please get the letters to me soon
- Join piazza: https://piazza.com/umd/fall2025/cmsc416cmsc616
- Assignment 0.1 will be posted tonight Sep 11 11:59 pm, due on Sep 18 11:59 pm

Writing parallel programs

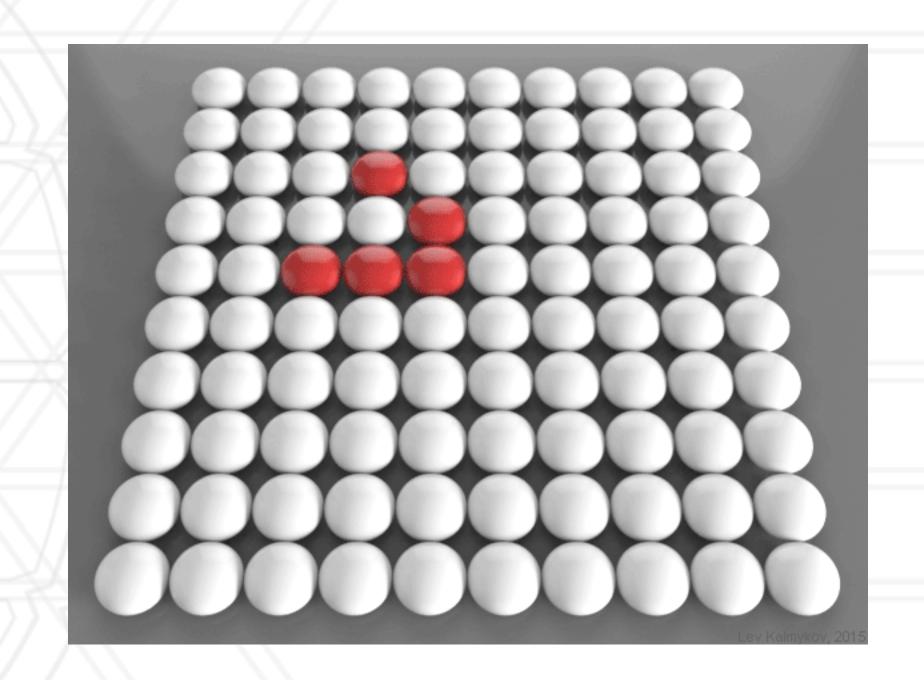
SPMD model

- Decide the serial algorithm first
- Data: how to distribute data among threads/processes?
 - Data locality: assignment of data to specific processes to minimize data movement
- Computation: how to divide work among threads/processes?
- Figure out how often communication will be needed



Conway's Game of Life

- Two-dimensional grid of (square) cells
- Each cell can be in one of two states: live or dead
- Every cell only interacts with its eight nearest neighbors
- In every generation (or iteration or time step), there are some rules that decide if a cell will continue to live or die or be born (dead → live)



https://en.wikipedia.org/wiki/Conway's_Game_of_Life

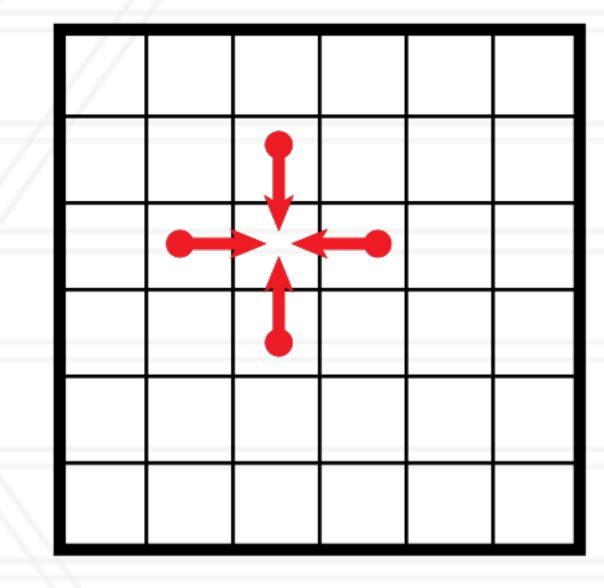
By Lev Kalmykov - Own work, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=43448735

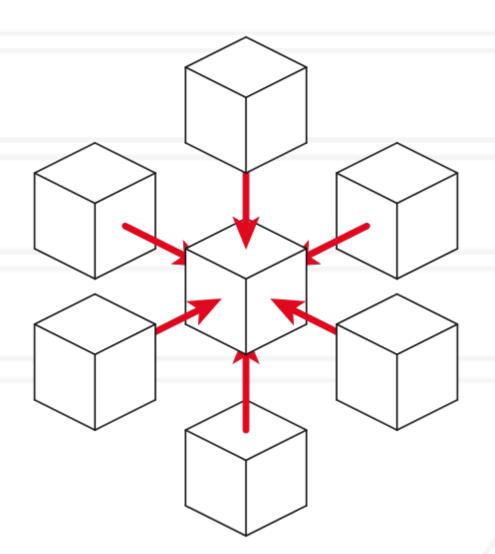


Two-dimensional stencil computation

2D 5-point Stencil

- Commonly found kernel in computational codes
- Heat diffusion, Jacobi method, Gauss-Seidel method





$$A[i,j] = \frac{A[i,j] + A[i-1,j] + A[i+1,j] + A[i,j-1] + A[i,j+1]}{5}$$

3D 7-point Stencil



Serial code

```
for(int t=0; t<num_steps; t++) {
    ...

// copy contents of A_new into A

for(i ...)
    for(j ...)
    A_new[i, j] = (A[i, j] + A[i-1, j] + A[i+1, j] + A[i, j-1] + A[i, j+1]) * 0.2</pre>
```

}

For correctness, we have to ensure that elements in A are not written into before they are read in the same timestep / iteration

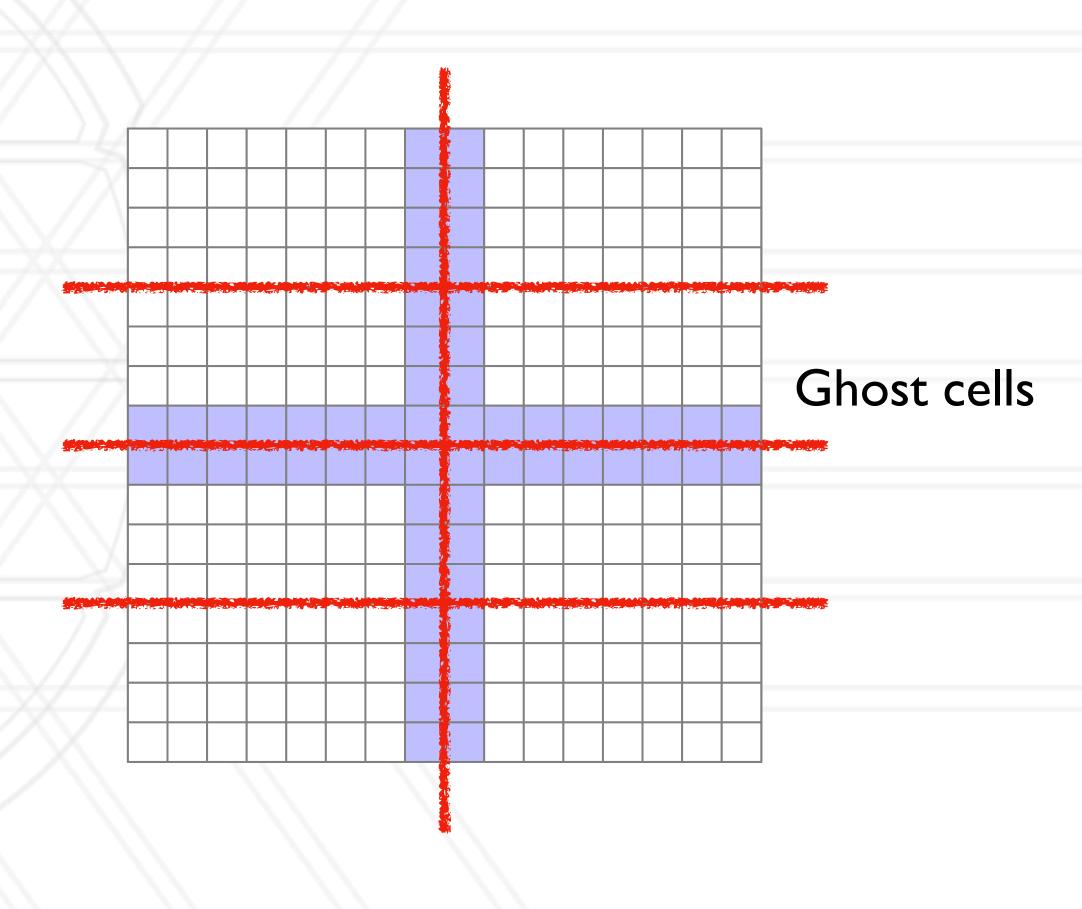


Why do we keep two

2D stencil computation in parallel

ID decomposition

- Divide rows (or columns) among processes
- Each process has to communicate with two neighbors (above and below)
- 2D decomposition
 - Divide both rows and columns (2d blocks) among processes
 - Each process has to communicate with four neighbors





Announcements

- Office hours are posted on the website
- Reminder: Assignment 0.1 is due on: Sep 18, 11:59 pm ET
- Good-faith attempt of each assignment is required
- Reminders:
 - How to contact course staff: cmsc416@cs.umd.edu
 - Mention your course and section number
 - Do not run/execute code on the login node
 - Do not run sudo on zaratan
 - Best way to report issues: What command did you run, and the full error



Prefix sum

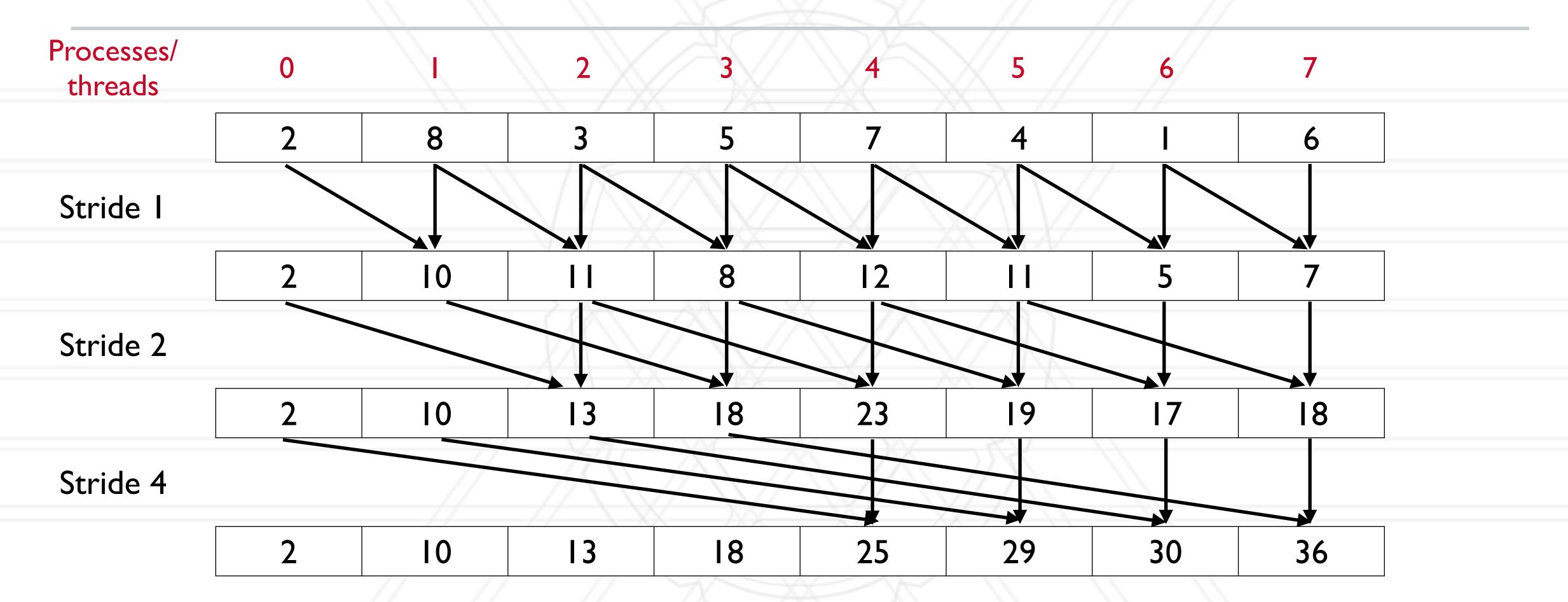
- Calculate sums of prefixes (running totals) of elements (numbers) in an array
- Also called a "scan" sometimes

```
pSum[0] = A[0]

for(i=1; i<N; i++) {
    pSum[i] = pSum[i-1] + A[i]
}</pre>
```

```
A 1 2 3 4 5 6 ...
pSum 1 3 6 10 15 21 ...
```

Parallel prefix sum





In practice

- You have N numbers and p processes, N >> p
- Assign a N/p block to each process
 - Do the serial prefix sum calculation for the blocks owned on each process locally
- Then do parallel algorithm with partial prefix sums (using the last element from each local block)
 - Last element from sending process is added to all elements in receiving process' sub-block

Load balance and grain size

- Load balance: try to balance the amount of work (computation) assigned to different threads/ processes
 - Bring ratio of maximum to average load as close to 1.0 as possible
 - Secondary consideration: also load balance amount of communication
- Grain size: ratio of computation-to-communication
 - Coarse-grained (more computation) vs. fine-grained (more communication)





Abhinav Bhatele

5218 Brendan Iribe Center (IRB) / College Park, MD 20742

phone: 301.405.4507 / e-mail: bhatele@cs.umd.edu