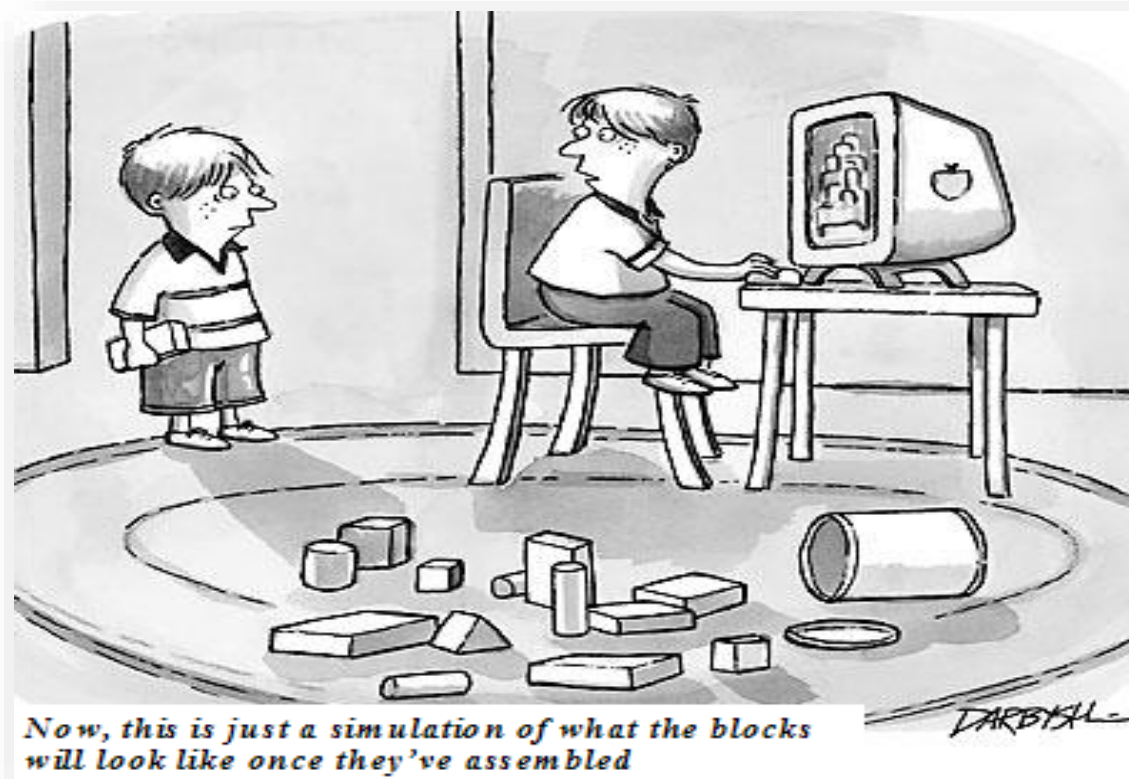


ECE569

Module 48



- MPI

Coding Styles

- **SPMD (Single Program, Multiple Data)**

- All PE's (Processor Elements) execute the same program in parallel, but has its own data
- CUDA Grid model (also OpenCL, MPI)
- SIMD is a special case – WARP used for efficiency

- **Master/Worker**

- A Master thread sets up a pool of worker threads
- Workers execute concurrently, removing tasks until done

- **Loop Parallelism (OpenMP)**

- Loop iterations execute in parallel
- FORTRAN do-all (truly parallel), do-across (with dependence)

- **Fork/Join**

- Most general, generic way of creation of threads

Program Models

SPMD

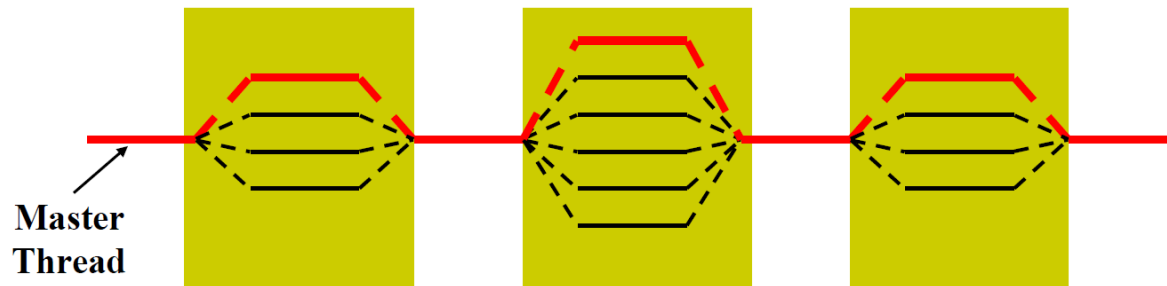
Master/Worker

Loop Parallelism

Fork/Join

OpenMP Programming Model

- Master thread spawns a team of threads as needed
 - Managed transparently on your behalf
 - relies on low-level thread fork/join methodology to implement parallelism
 - The developer is spared the details
- Leveraging OpenMP in an existing code:
 - Parallelism is added incrementally: that is, the sequential program evolves into a parallel program
- **Not Scalable**



SPMD Program

- **Dominant coding style of scalable computing**
 - MPI code is mostly developed in SPMD style
 - Many OpenMP code is also in SPMD
 - Particularly suitable for algorithms based on task parallelism and geometric decomposition.

SPMD is by far the most commonly used pattern for structuring massively parallel programs.

SPMD Program Phases

- **Initialize**
 - Establish localized data structure and communication
- **Obtain a unique identifier**
 - Each thread acquires a unique identifier, typically range from 0 to $N-1$, where N is the number of threads.
 - Both OpenMP and CUDA have built-in support for this.
- **Distribute Data**
 - Decompose global data into chunks and localize them, or
 - Sharing/replicating major data structure using thread ID to associate subset of the data to threads
- **Run the core computation**
- **Finalize**
 - Reconcile global data structure,

Data Sharing

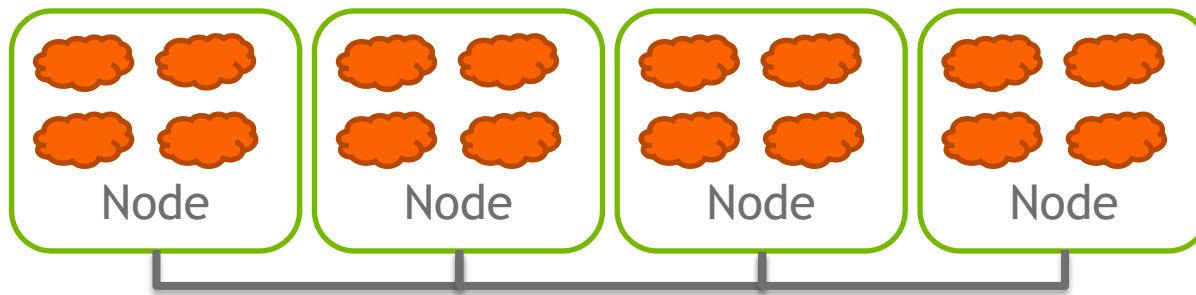
- **Data sharing can be a double-edged sword**
 - Excessive data sharing drastically reduces advantage of parallel execution
 - Localized sharing can improve memory bandwidth efficiency
- **Efficient memory bandwidth usage can be achieved by synchronizing the execution of task groups and coordinating their usage of memory data**
 - Efficient use of on-chip, shared storage and datapaths
- **Read-only sharing can usually be done at much higher efficiency than read-write sharing, which often requires more synchronization**

MPI

- A distributed memory model
 - processes exchange information by messaging
- API communication functions
 - Seamless interconnect network
- Processes address each other using logical numbers
- No cache coherence and no need for special cache coherency hardware
- Software development: more difficult to write programs: keep track of memory usage

MPI – Execution Model

- Many processes distributed in a cluster
- Each process computes part of the output
- Processes communicate with each other
- Processes can synchronize when collaborating on a large task
 - What differentiates processes is their rank: “branching based on the process rank”
- Very similar to GPU computing, where one thread did work based on its thread index

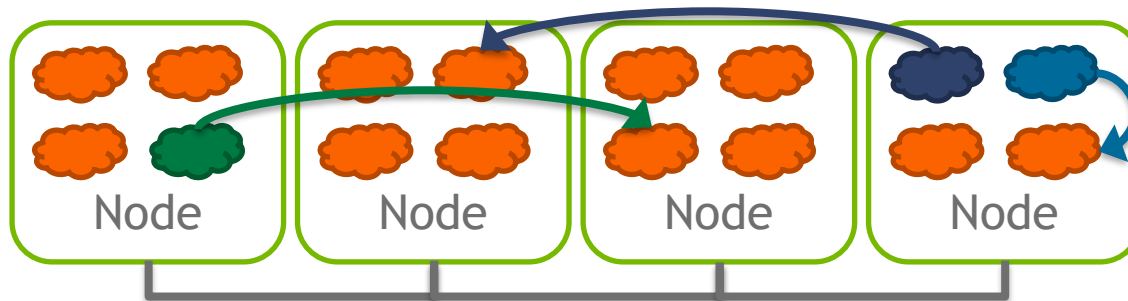


MPI Initialization, Info and Sync

- `int MPI_Init(int *argc, char ***argv)`
 - Initialize MPI
- `int MPI_Comm_rank (MPI_Comm comm, int *rank)`
 - Rank of the calling process in group of comm
- `int MPI_Comm_size (MPI_Comm comm, int *size)`
 - Number of processes in the group of comm
- `int MPI_Abort (MPI_Comm comm)`
 - Terminate MPI communication with an error flag
- `int MPI_Finalize ()`
 - Ending an MPI application, close all resources

MPI Point to Point Communication: Sending Data

- **int MPI_Send(void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)**
 - Buf: Initial address of send buffer (choice)
 - Count: Number of elements in send buffer (nonnegative integer)
 - Datatype: Datatype of each send buffer element (handle)
 - Dest: Rank of destination (integer)
 - Tag: Message tag (integer)
 - Comm: Communicator (handle)



MPI Point to Point Communication: Receiving Data

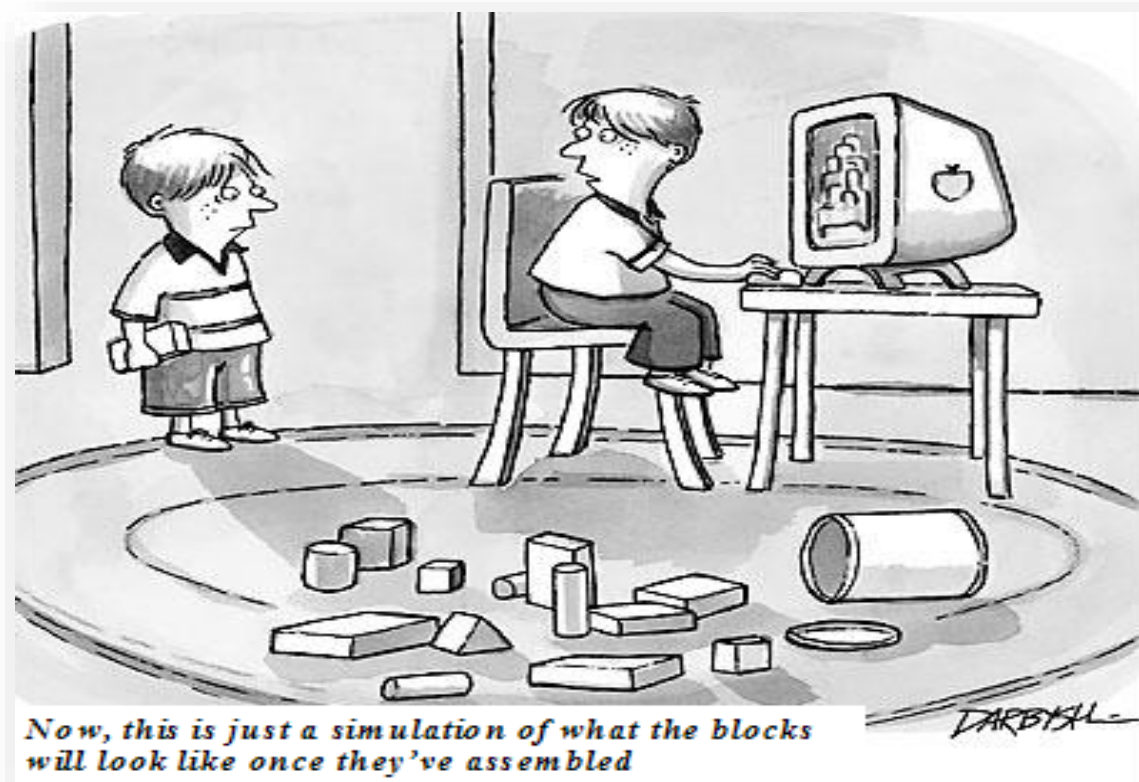
- `int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status *status)`
 - Buf: Initial address of receive buffer (choice)
 - Count: Maximum number of elements in receive buffer (integer)
 - Datatype: Datatype of each receive buffer element (handle)
 - Source: Rank of source (integer)
 - Tag: Message tag (integer) – `MPI_ANY_TAG`
 - Comm: Communicator (handle)
 - Status: Status object (Status)

Next

- **MPI Example**

ECE569

Module 49



- MPI Example: Vector Addition

MPI Example

```
1  #include "mpi.h"
2  #include <stdio.h>
3  #include <string.h>
4  int main(int argc, char* argv[]) {
5  int my_rank; /* rank of process */
6  int p; /* number of processes */
7  int source; /* rank of sender */
8  int dest; /* rank of receiver */
9  int tag = 0; /* tag for messages */
10 char message[100]; /* storage for message */
11 MPI_Status status; /* return status for receive */

12 MPI_Init(&argc, &argv); // Start up MPI
13 MPI_Comm_rank(MPI_COMM_WORLD, &my_rank); // Find out process rank
14 MPI_Comm_size(MPI_COMM_WORLD, &p); // Find out number of processes
```

MPI Example

```
15  if (my_rank != 0) {
16      /* Create message */
17      sprintf(message, "Greetings from process %d!", my_rank);
18      dest = 0;
19      /* Use strlen+1 so that '\0' gets transmitted */
21      MPI_Send(message, strlen(message)+1, MPI_CHAR, dest, tag,
                                                         MPI_COMM_WORLD);
22  }
23  else { /* my_rank == 0 */
24      for (source = 1; source < p; source++) {
25          MPI_Recv(message, 100, MPI_CHAR, source, tag, MPI_COMM_WORLD,
                  &status);
26          printf("%s\n", message);
27      }
28  }
29  MPI_Finalize(); // Shut down MPI
30  return 0;
31 } /* main */
```

MPI Example – Vector Addition (main)

```
1  #include "mpi.h"
2  #include <stdio.h>
3  #include <string.h>
4  int main(int argc, char* argv[]) {
5  int vector_size=1024*1024*1024;
6  int np; /* number of processes */
7  int pid; /* rank */

8  MPI_Init(&argc, &argv); // Start up MPI
9  MPI_Comm_rank(MPI_COMM_WORLD, &pid); // Find out process rank
10 MPI_Comm_size(MPI_COMM_WORLD, &np); // Find out number of processes

11 If (np < 3 ) {
12     if (pid == 0) printf ("need 3 or more processes\n");
13     MPI_Abort(MPI_COMM_WORLD,1);
14     return 1;
15 }
```


MPI Example – Vector Addition (main)

```
16 if (pid < np-1)
17     compute_node(vector_size/(np-1));
18 else
19     data_server(vector_size)
20 MPI_Finalize(); // Shut down MPI
21 return 0;
22 } /* main */
```

MPI Example Vector Addition: Server(1)

```
void data_server(unsigned int vector_size) {
1   int np, num_nodes = np - 1;
2   unsigned int num_bytes = vector_size * sizeof(float);
3   float *input_a = 0, *input_b = 0, *output = 0;

4   /* Set MPI Communication Size */
5   MPI_Comm_size(MPI_COMM_WORLD, &np);

6   /* Allocate input data */
7   input_a = (float *)malloc(num_bytes);
8   input_b = (float *)malloc(num_bytes);
9   output = (float *)malloc(num_bytes);
10  if(input_a == NULL || input_b == NULL || output == NULL) {
11      printf("Server couldn't allocate memory\n");
12      MPI_Abort( MPI_COMM_WORLD, 1 );
13  }
14  /* Initialize input data */
15  random_data(input_a, vector_size , 1, 10);
16  random_data(input_b, vector_size , 1, 10);
```

MPI Example Vector Addition: Server(2)

```
17  /* Send data to compute nodes */
18  float *ptr_a = input_a;
19  float *ptr_b = input_b;

20  for(int process=0; process < num_nodes; process++) {
21      MPI_Send(ptr_a, vector_size / num_nodes, MPI_FLOAT,
22              process, DATA_DISTRIBUTE, MPI_COMM_WORLD);
23      ptr_a += vector_size / num_nodes;

24      MPI_Send(ptr_b, vector_size / num_nodes, MPI_FLOAT,
25              process, DATA_DISTRIBUTE, MPI_COMM_WORLD);
26      ptr_b += vector_size / num_nodes;
27  }

28  /* Wait for nodes to compute */
29  MPI_Barrier(MPI_COMM_WORLD);
```

MPI Example Vector Addition: Server(3)

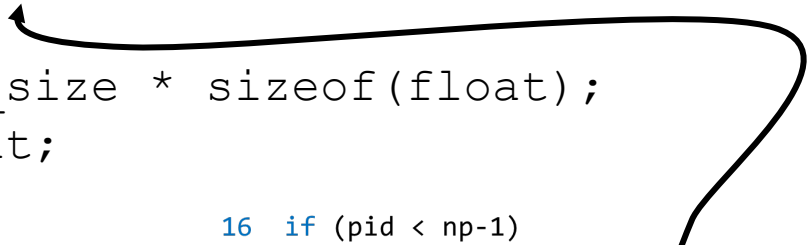
```
28  /* Serve node ready to receive data from compute nodes */
29  /* set up MPI_Status flag for error check
30  /* Collect output data */
31  MPI_Status status;
32  for(int process = 0; process < num_nodes; process++) {
33      MPI_Recv(output + process * vector_size / num_nodes,
               vector_size / num_comp_nodes, MPI_FLOAT, process,
               DATA_COLLECT, MPI_COMM_WORLD, &status );
34  }

35  /* Store output data */
36  print_output(output, vector_size);

37  /* Release resources */
38  free(input_a);
39  free(input_b);
40  free(output);
41 }
```

MPI Example Vector Addition: Compute(1)

```
void compute_node(unsigned int vector_size ) {  
1  int np;  
2  unsigned int num_bytes = vector_size * sizeof(float);  
3  float *input_a, *input_b, *output;  
4  MPI_Status status;  
  
5  MPI_Comm_size(MPI_COMM_WORLD, &np);  
6  int server_process = np - 1;  
  
7  /* Alloc host memory */  
8  input_a = (float *)malloc(num_bytes);  
9  input_b = (float *)malloc(num_bytes);  
10 output = (float *)malloc(num_bytes);  
  
11 /* Get the input data from server process */  
12 MPI_Recv(input_a, vector_size, MPI_FLOAT,  
           server_process, DATA_DISTRIBUTE, MPI_COMM_WORLD,  
           &status);  
13 MPI_Recv(input_b, vector_size, MPI_FLOAT,  
           server_process, DATA_DISTRIBUTE, MPI_COMM_WORLD,  
           &status);  
  
16 if (pid < np-1)  
17     compute_node(vector_size/(np-1));  
18 else  
19     data_server(vector_size)  
20 MPI_Finalize(); // Shut down MPI  
21 return 0;  
22 } /* main */
```



MPI Example Vector Addition: Compute(2)

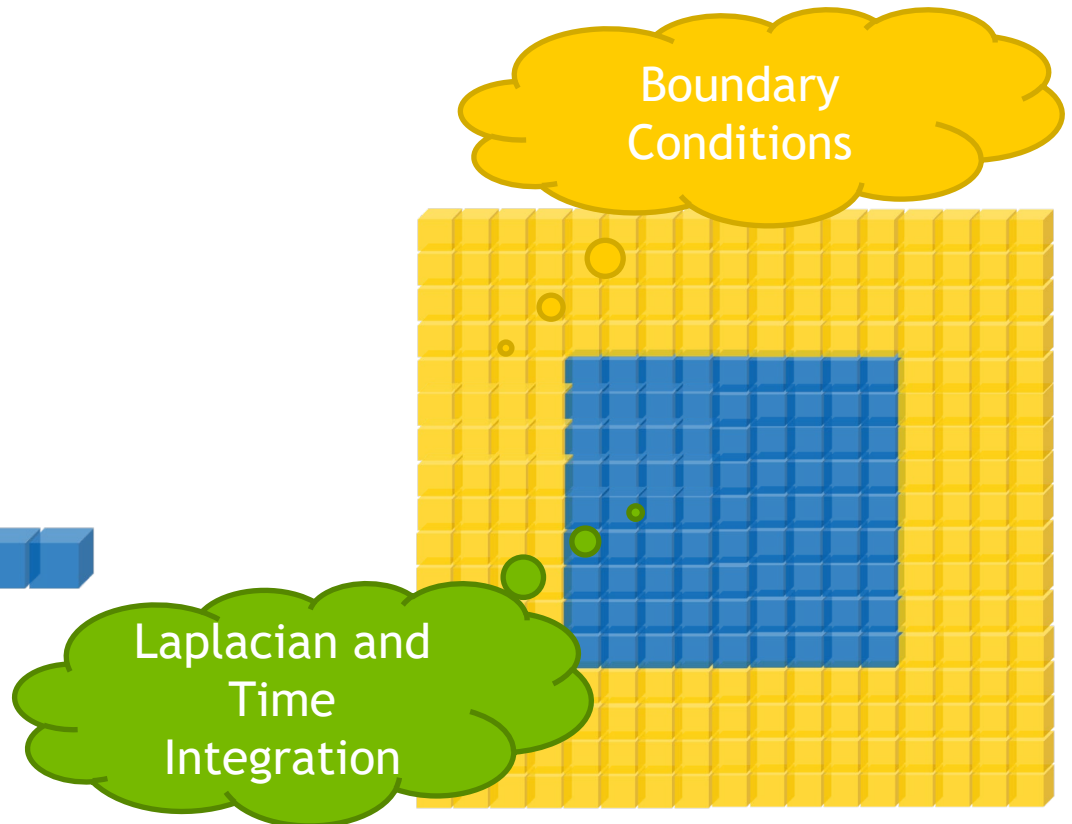
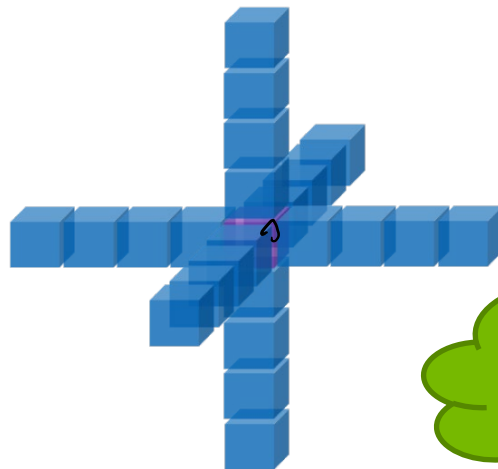
```
14  /* Compute the partial vector addition */
15  for(int i = 0; i < vector_size; ++i) {
16      output[i] = input_a[i] + input_b[i];
17  }
18  /* Replace serial execution with CUDA kernel */
19  /* same as your vector addition main function */
20  /* and kernel code */
21  /* Report to barrier after computation is done*/
22  MPI_Barrier(MPI_COMM_WORLD);

23  /* Send the output */
24  MPI_Send(output, vector_size, MPI_FLOAT,
           server_process, DATA_COLLECT, MPI_COMM_WORLD);

25  /* Release memory */
26  free(input_a);
27  free(input_b);
28  free(output);
29 }
```

Next: Wave Propagation Example

- **Approximate Laplacian using finite differences**
- 3D Stencil, 4 points in each direction (x,y,z)
 - $(i-4,j,k), (i-3,j,k), (i-2,j,k), (i-1,j,k), (i+1,j,k), (i+2,j,k), (i+3,j,k), (i+4,j,k), (i,j-4,k), (i,j-3,k), (i,j-2,k), (i,j-1,k), (i,j+1,k), (i,j+2,k), (i,j+3,k), (i,j+4,k), (i,j,k-4), (i,j,k-3), (i,j,k-2), (i,j,k-1), (i,j,k+1), (i,j,k+2), (i,j,k+3)$ and $(i,j,k+4)$.



Next: Stencil Domain Decomposition

- **Volumes are split into tiles (along the Z-axis)**
 - 3D-Stencil introduces data dependencies

