

CHAPTER 20

Programming a heterogeneous computing cluster

An introduction to CUDA streams

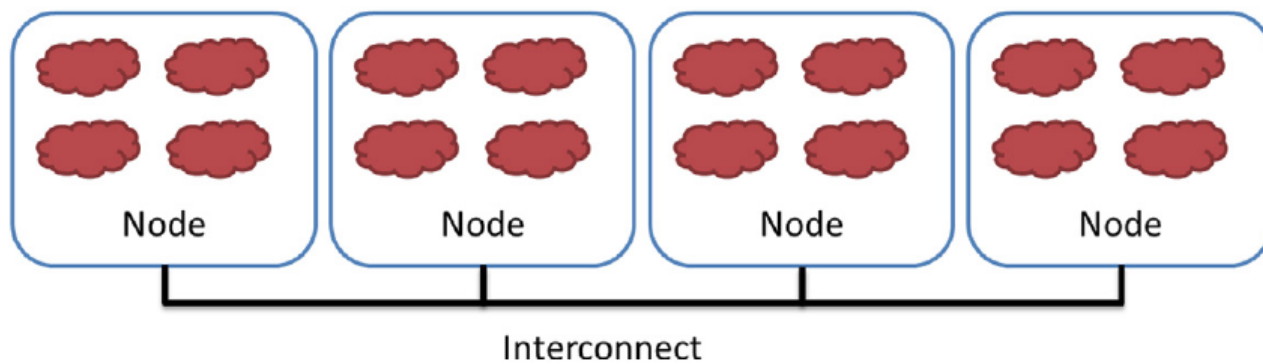
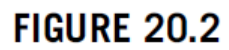


FIGURE 20.1

Programmer's view of MPI processes. *MPI*, Message Passing Interface.



A 25-point stencil computation example, with four neighbors in each of the x , y , and z directions.

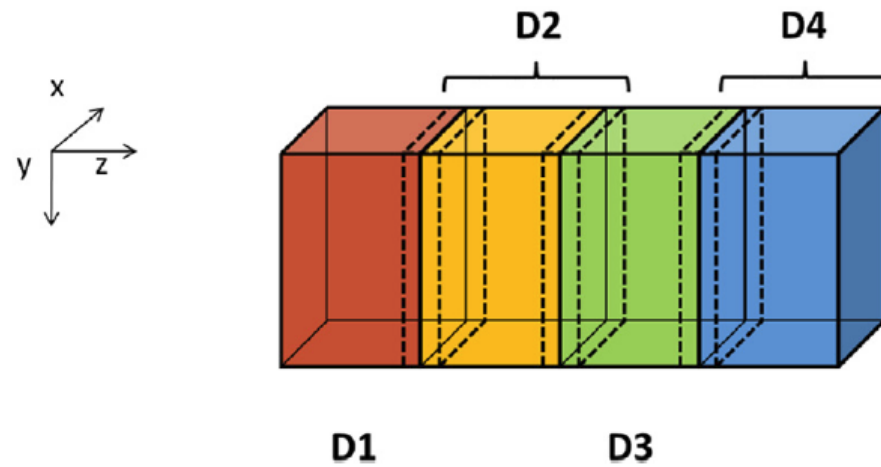


FIGURE 20.3

3D Grid array for the modeling heat transfer in a duct.

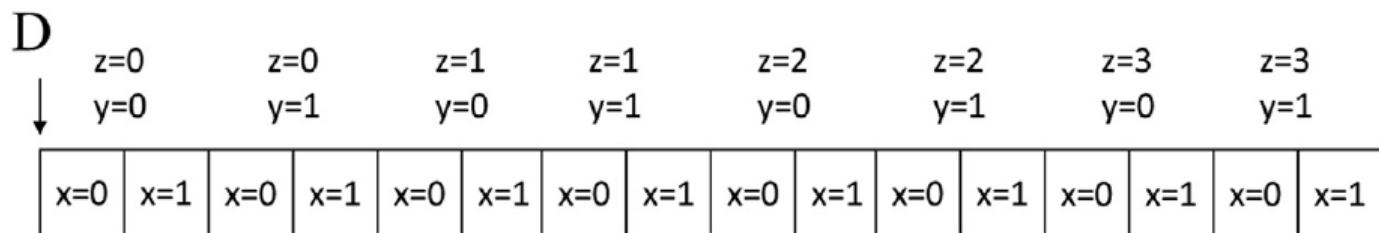


FIGURE 20.4

A small example of memory layout for the 3D grid.

- `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_Comm_abort (MPI_Comm comm)`
 - Terminate MPI communication connection with an error flag
- `int MPI_Finalize ()`
 - Ending an MPI application, close all resources

FIGURE 20.5

Basic MPI functions for establishing and closing a communication system.

```

01  #include "mpi.h"
02  int main(int argc, char *argv[]) {
03      int pad = 0, dimx = 480+pad, dimy = 480, dimz = 400, nreps = 100;
04      int pid=-1, np=-1;
05      MPI_Init(&argc, &argv);
06      MPI_Comm_rank(MPI_COMM_WORLD, &pid);
07      MPI_Comm_size(MPI_COMM_WORLD, &np);
08      if(np < 3) {
09          if(0 == pid) printf("Needed 3 or more processes.\n");
10          MPI_Abort( MPI_COMM_WORLD, 1 ); return 1;
11      }
12      if(pid < np - 1)
13          compute_process(dimx, dimy, dimz / (np - 1), nreps);
14      else
15          data_server( dimx,dimy,dimz );
16      MPI_Finalize();
17      return 0;
18  }

```

FIGURE 20.6

A simple MPI main program.

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

FIGURE 20.7

Syntax for the MPI_Send() function.

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

FIGURE 20.8

Syntax for the MPI_Recv() function.

```

01 void data_server(int dimx, int dimy, int dimz, int nreps) {
02     int np;
03     /* Set MPI Communication Size */
04     MPI_Comm_size(MPI_COMM_WORLD, &np);
05     unsigned int num_comp_nodes = np - 1, first_node = 0, last_node = np - 2;
06     unsigned int num_points = dimx * dimy * dimz;
07     unsigned int num_bytes = num_points * sizeof(float);
08     float *input=0, *output=0;
09     /* Allocate input data */
10     input = (float *)malloc(num_bytes);
11     output = (float *)malloc(num_bytes);
12     if(input == NULL || output == NULL) {
13         printf("server couldn't allocate memory\n");
14         MPI_Abort( MPI_COMM_WORLD, 1 );
15     }
16     /* Initialize input data */
17     random_data(input, dimx, dimy, dimz, 1, 10);
18     /* Calculate number of shared points */
19     int edge_num_points = dimx * dimy * ((dimz / num_comp_nodes) + 4);
20     int int_num_points = dimx * dimy * ((dimz / num_comp_nodes) + 8);
21     float *send_address = input;
22     /* Send data to the first compute node */
23     MPI_Send(send_address, edge_num_points, MPI_FLOAT, first_node,
24             0, MPI_COMM_WORLD);
25     send_address += dimx * dimy * ((dimz / num_comp_nodes) - 4);
26     /* Send data to "internal" compute nodes */
27     for(int process = 1; process < last_node; process++) {
28         MPI_Send(send_address, int_num_points, MPI_FLOAT, process,
29             0, MPI_COMM_WORLD);
30         send_address += dimx * dimy * (dimz / num_comp_nodes);
31     }
32     /* Send data to the last compute node */
33     MPI_Send(send_address, edge_num_points, MPI_FLOAT, last_node,
34             0, MPI_COMM_WORLD);

```

FIGURE 20.9

Data server process code (part 1).

```

01 void compute_node_stencil(int dimx, int dimy, int dimz, int nreps ) {
02     int np, pid;
03     MPI_Comm_rank(MPI_COMM_WORLD, &pid);
04     MPI_Comm_size(MPI_COMM_WORLD, &np);
05     int server_process = np - 1;
06     unsigned int num_points      = dimx * dimy * (dimz + 8);
07     unsigned int num_bytes      = num_points * sizeof(float);
08     unsigned int num_halo_points = 4 * dimx * dimy;
09     unsigned int num_halo_bytes = num_halo_points * sizeof(float);
10     /* Allocate host memory */
11     float *h_input = (float *)malloc(num_bytes);
12     /* Allocate device memory for input and output data */
13     float *d_input = NULL;
14     cudaMalloc((void **)&d_input, num_bytes );
15     float *rcv_address = h_input + ((0 == pid) ? num_halo_points : 0);
16     MPI_Recv(rcv_address, num_points, MPI_FLOAT, server_process,
17             MPI_ANY_TAG, MPI_COMM_WORLD, &status );
18     cudaMemcpy(d_input, h_input, num_bytes, cudaMemcpyHostToDevice);

```

FIGURE 20.10

Compute process code (part 1).

```

16     float *h_output = NULL, *d_output = NULL, *d_vsq = NULL;
17     float *h_output = (float *)malloc(num_bytes);
18     cudaMalloc((void **)&d_output, num_bytes );
19     float *h_left_boundary = NULL, *h_right_boundary = NULL;
20     float *h_left_halo = NULL, *h_right_halo = NULL;
21     /* Allocate host memory for halo data */
22     cudaHostAlloc((void **)&h_left_boundary, num_halo_bytes,
23                  cudaHostAllocDefault);
24     cudaHostAlloc((void **)&h_right_boundary, num_halo_bytes,
25                  cudaHostAllocDefault);
26     cudaHostAlloc((void **)&h_left_halo, num_halo_bytes,
27                  cudaHostAllocDefault);
28     cudaHostAlloc((void **)&h_right_halo, num_halo_bytes,
29                  cudaHostAllocDefault);
30     /* Create streams used for stencil computation */
31     cudaStream_t stream0, stream1;
32     cudaStreamCreate(&stream0);
33     cudaStreamCreate(&stream1);

```

FIGURE 20.11

Compute process code (part 2).

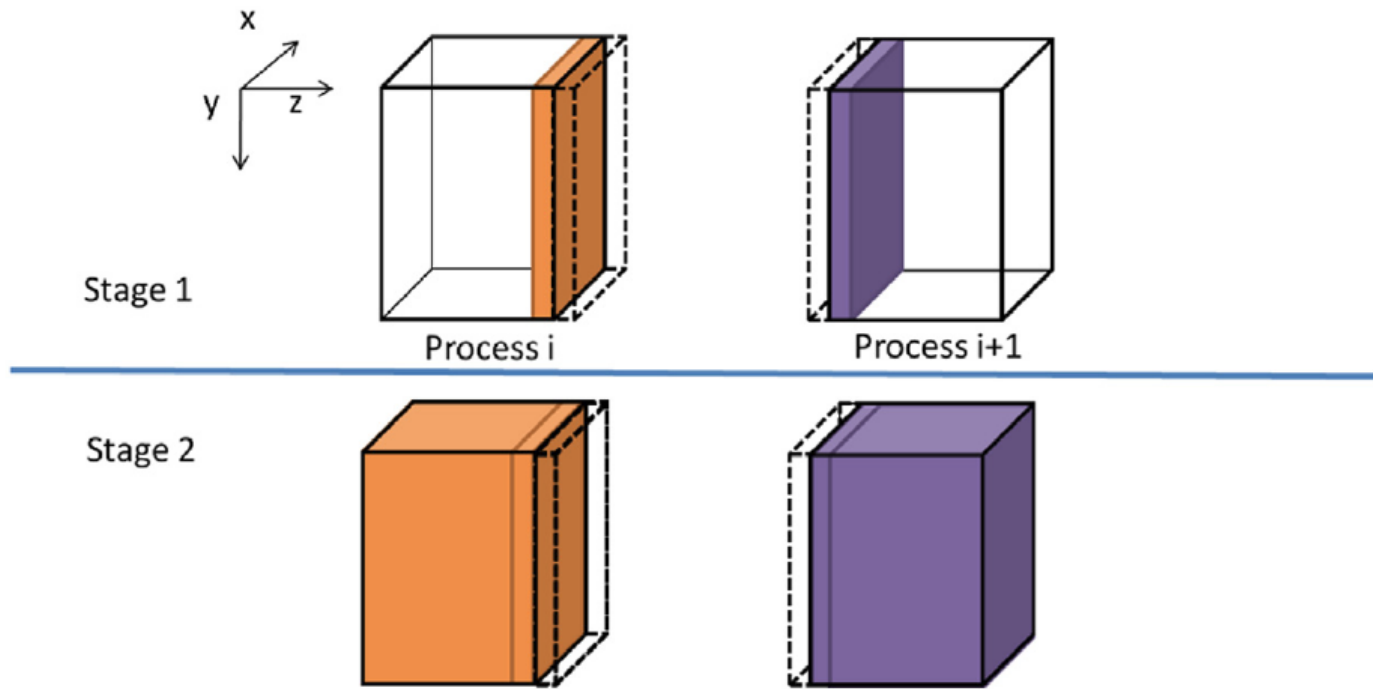


FIGURE 20.12

A two-stage strategy for overlapping computation with communication.

```

28 MPI_Status status;
29 int left_neighbor = (pid > 0)      ? (pid - 1) : MPI_PROC_NULL;
30 int right_neighbor = (pid < np - 2) ? (pid + 1) : MPI_PROC_NULL;

    /* Upload stencil coefficients */
31 upload_coefficients(coeff, 5);
32 int left_halo_offset = 0;
33 int right_halo_offset = dimx * dimy * (4 + dimz);
34 int left_stag1_offset = 0;
35 int right_stag1_offset = dimx * dimy * (dimz - 4);
36 int stage2_offset = num_halo_points;
37 MPI_Barrier( MPI_COMM_WORLD );
38 for(int i=0; i < nreps; i++) {
    /* Compute boundary values needed by other nodes first */
39     call_stencil_kernel(d_output + left_stag1_offset,
        d_input + left_stag1_offset, dimx, dimy, 12, stream0);
40     call_stencil_kernel(d_output + right_stag1_offset,
        d_input + right_stag1_offset, dimx, dimy, 12, stream0);
    /* Compute the remaining points */
41     call_stencil_kernel(d_output + stage2_offset, d_input +
        stage2_offset, dimx, dimy, dimz, stream1);

```

FIGURE 20.13

Compute process code (part 3).

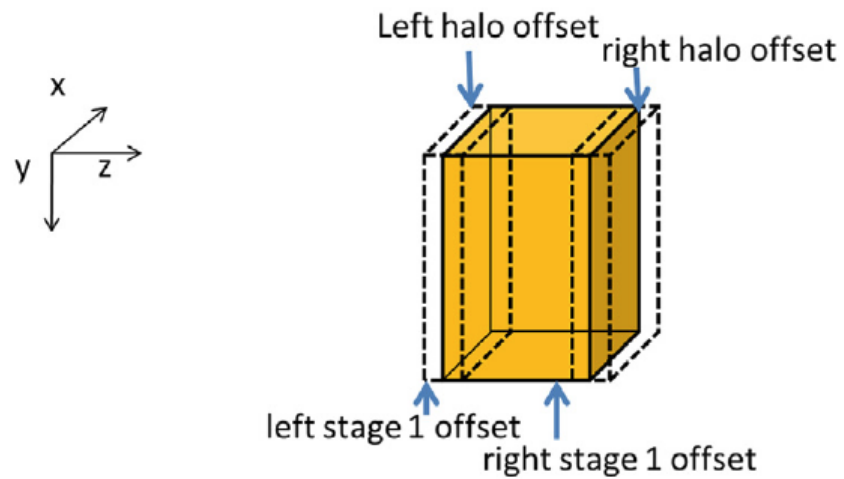


FIGURE 20.14

Device memory offsets used for data exchange with neighbor processes.

```

42      /* Copy the data needed by other nodes to the host */
      cudaMemcpyAsync(h_left_boundary, d_output + num_halo_points,
                     num_halo_bytes, cudaMemcpyDeviceToHost, stream0 );
43      cudaMemcpyAsync(h_right_boundary,
                     d_output + right_stagel_offset + num_halo_points,
                     num_halo_bytes, cudaMemcpyDeviceToHost, stream0 );
44      cudaStreamSynchronize(stream0);
      /* Send data to left, get data from right */
45      MPI_Sendrecv(h_left_boundary, num_halo_points, MPI_FLOAT,
                   left_neighbor, i, h_right_halo, num_halo_points,
                   MPI_FLOAT, right_neighbor, i, MPI_COMM_WORLD, &status );
      /* Send data to right, get data from left */
46      MPI_Sendrecv(h_right_boundary, num_halo_points, MPI_FLOAT,
                   right_neighbor, i, h_left_halo, num_halo_points,
                   MPI_FLOAT, left_neighbor, i, MPI_COMM_WORLD, &status );

47      cudaMemcpyAsync(d_output+left_halo_offset, h_left_halo,
                     num_halo_bytes, cudaMemcpyHostToDevice, stream0);
48      cudaMemcpyAsync(d_output+right_halo_offset, h_right_halo,
                     num_halo_bytes, cudaMemcpyHostToDevice, stream0 );
49      cudaDeviceSynchronize();

50      float *temp = d_output;
51      d_output = d_input; d_input = temp;
52  }

```

FIGURE 20.15

Compute process code (part 4).

- `int MPI_Sendrecv(void *sendbuf, int sendcount, MPI_Datatype sendtype, int dest, int sendtag, void *recvbuf, int recvcount, MPI_Datatype recvtype, int source, int recvtag, MPI_Comm comm, MPI_Status *status)`
 - **Sendbuf**: Initial address of send buffer (choice)
 - **Sendcount**: Number of elements in send buffer (integer)
 - **Sendtype**: Type of elements in send buffer (handle)
 - **Dest**: Rank of destination (integer)
 - **Sendtag**: Send tag (integer)
 - **Recvcount**: Number of elements in receive buffer (integer)
 - **Recvtype**: Type of elements in receive buffer (handle)
 - **Source**: Rank of source (integer)
 - **Recvtag**: Receive tag (integer)
 - **Comm**: Communicator (handle)
 - **Recvbuf**: Initial address of receive buffer (choice)
 - **Status**: Status object (Status). This refers to the receive operation.

FIGURE 20.16

Syntax for the `MPI_Sendrecv()` function.

```

    /* Wait for previous communications */
53  MPI_Barrier(MPI_COMM_WORLD);

54  float *temp = d_output;
55  d_output = d_input;
56  d_input = temp;

    /* Send the output, skipping halo points */
57  cudaMemcpy(h_output, d_output, num_bytes, cudaMemcpyDeviceToHost);
    float *send_address = h_output + num_ghost_points;
58  MPI_Send(send_address, dimx * dimy * dimz, MPI_REAL,
            server_process, DATA_COLLECT, MPI_COMM_WORLD);
59  MPI_Barrier(MPI_COMM_WORLD);

    /* Release resources */
60  free(h_input); free(h_output);
61  cudaFreeHost(h_left_ghost_own); cudaFreeHost(h_right_ghost_own);
62  cudaFreeHost(h_left_ghost); cudaFreeHost(h_right_ghost);
63  cudaFree( d_input ); cudaFree( d_output );
64  }

```

FIGURE 20.17

Compute process code (part 5).

```
/* Wait for nodes to compute */
24 MPI_Barrier(MPI_COMM_WORLD);
/* Collect output data */
25 MPI_Status status;
26 for(int process = 0; process < num_comp_nodes; process++)
27     MPI_Recv(output + process * num_points / num_comp_nodes,
               num_points / num_comp_nodes, MPI_REAL, process,
               DATA_COLLECT, MPI_COMM_WORLD, &status );

/* Store output data */
28 store_output(output, dimx, dimy, dimz);
/* Release resources */
29 free(input);
30 free(output);
}
```

FIGURE 20.18

Data server code (part 2).

```
MPI_SendRecv(d_output + num_halo_points, num_halo_points, MPI_FLOAT,  
             left_neighbor, i, d_output + left_halo_offset, num_halo_points,  
             MPI_FLOAT, right_neighbor, i, MPI_COMM_WORLD, &status);  
MPI_SendRecv(d_output + right_stage1_offset, num_halo_points,  
             num_halo_points, MPI_FLOAT, right_neighbor, i,  
             d_output + right_halo_offset, num_halo_points,  
             MPI_FLOAT, left_neighbor, i, MPI_COMM_WORLD, &status);
```

FIGURE 20.19

Revised MPI SendRecv calls in using CUDA-aware MPI.