CS 7172 Parallel and Distributed Computation

Message Passing Interface (MPI)

Kun Suo

Computer Science, Kennesaw State University

https://kevinsuo.github.io/

Review

- MPI introduction
 - Helloworld of MPI

Performance evaluation

- Example: how to solve problems in MPI
 - Trapezoidal problem

Outline

MPI_Send/MPI_Receive

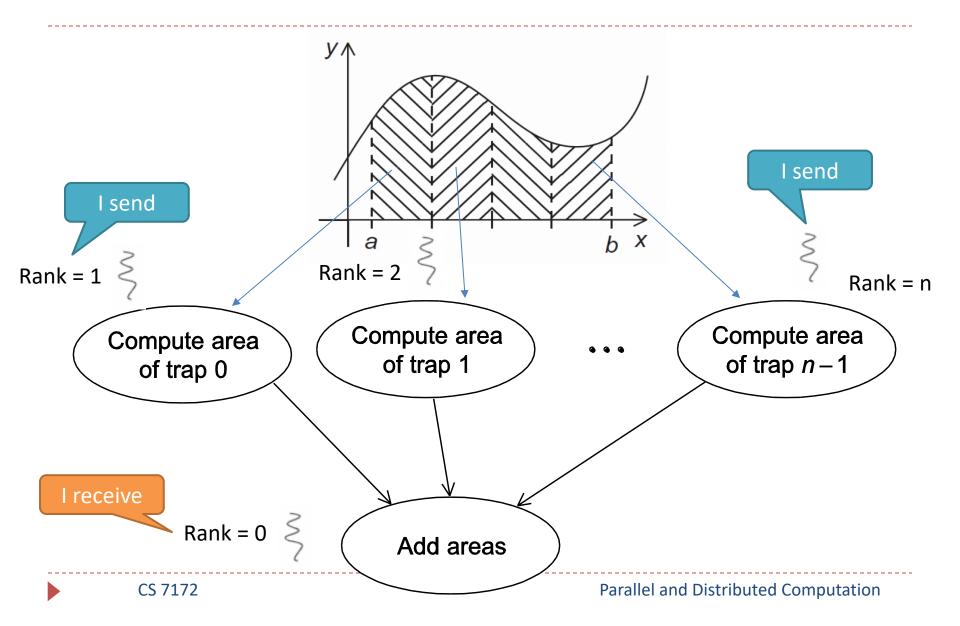
MPI_Reduce/MPI_Allreduce

MPI_Broadcast

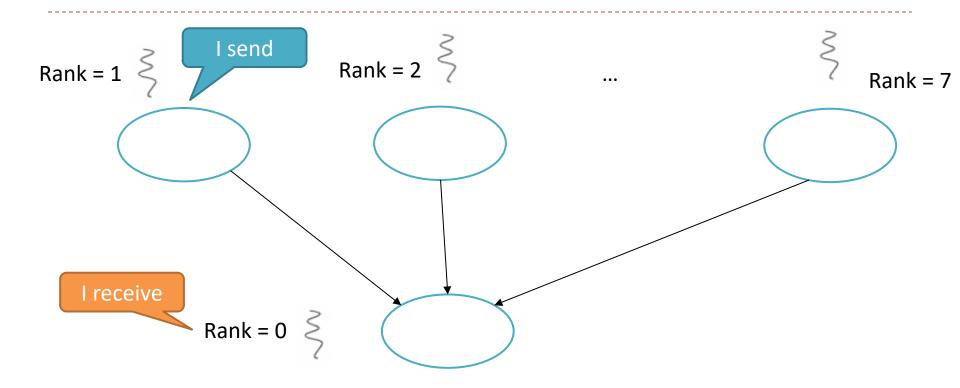
MPI_Scatter

MPI_Gather

Communication



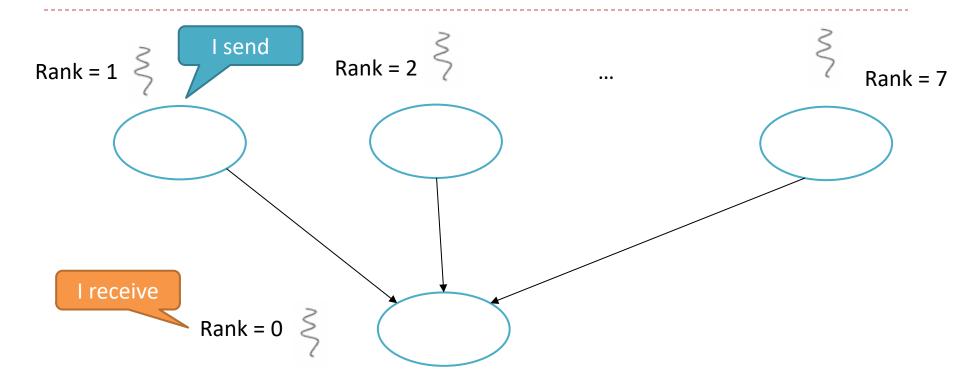
Point-to-Point communication



How many send/receive and add operation?

- 7 sends and 7 receives
- 7 adds

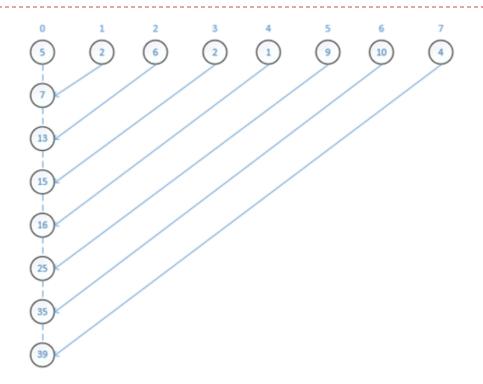
Point-to-Point communication



How many send/receive and add operation on master node (Rank = 0)?

- 7 sends and 7 receives
- 7 adds

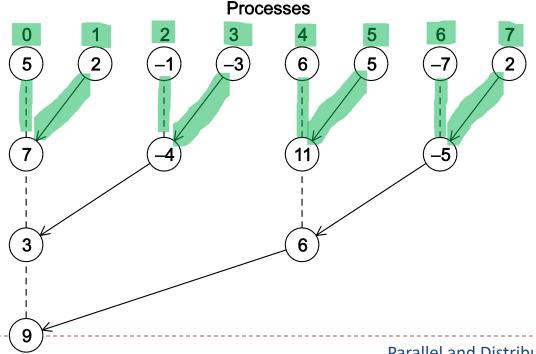
Point-to-Point communication



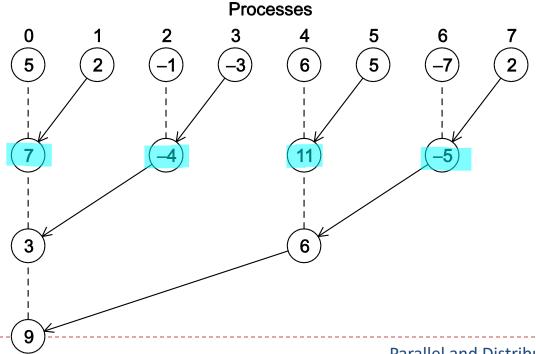
How many send/receive and add operation on master node (Rank = 0)?

- 7 sends and 7 receives
- 7 adds

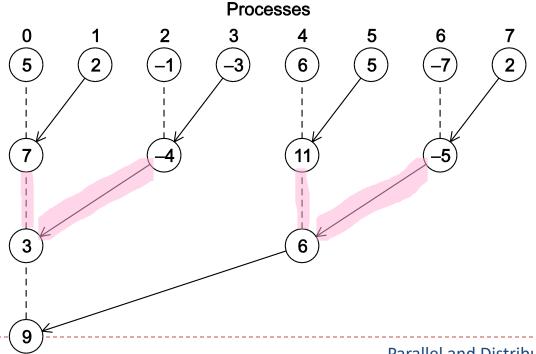
- o (a) Process 1 sends to 0; 3 sends to 2; 5 sends to 4; and 7 sends to 6.
- (b) Processes 0, 2, 4, and 6 add in the received values.
- (c) Processes 2 and 6 send their new values to processes 0 and 4, respectively.
- (d) Processes 0 and 4 add the received values into their new values.



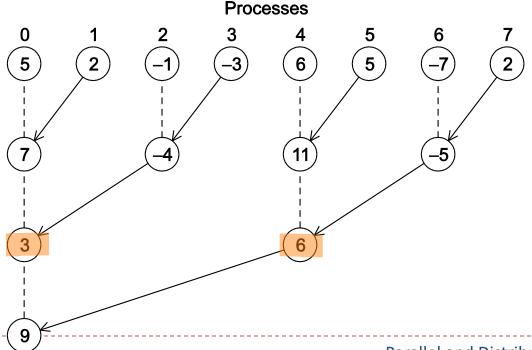
- o (a) Process 1 sends to 0; 3 sends to 2; 5 sends to 4; and 7 sends to 6.
- (b) Processes 0, 2, 4, and 6 add in the received values.
- (c) Processes 2 and 6 send their new values to processes 0 and 4, respectively.
- o (d) Processes 0 and 4 add the received values into their new values.



- (a) Process 1 sends to 0; 3 sends to 2; 5 sends to 4; and 7 sends to 6.
- o (b) Processes 0, 2, 4, and 6 add in the received values.
- (c) Processes 2 and 6 send their new values to processes 0 and 4, respectively.
- o (d) Processes 0 and 4 add the received values into their new values.

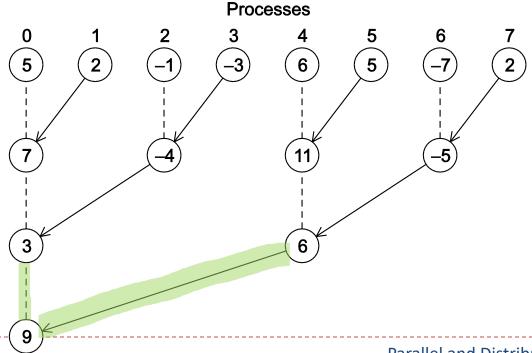


- (a) Process 1 sends to 0; 3 sends to 2; 5 sends to 4; and 7 sends to 6.
- (b) Processes 0, 2, 4, and 6 add in the received values.
- (c) Processes 2 and 6 send their new values to processes 0 and 4, respectively.
- o (d) Processes 0 and 4 add the received values into their new values.



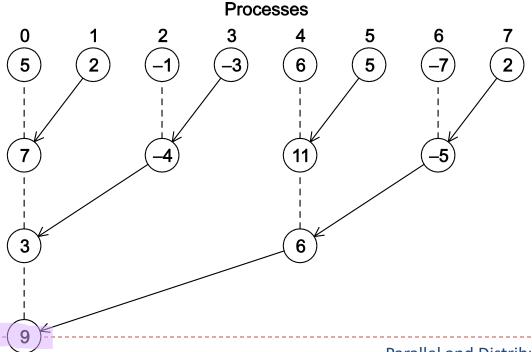
In the second phase:

- (a) Process 4 sends its newest value to process 0.
- (b) Process 0 adds the received value to its newest value.



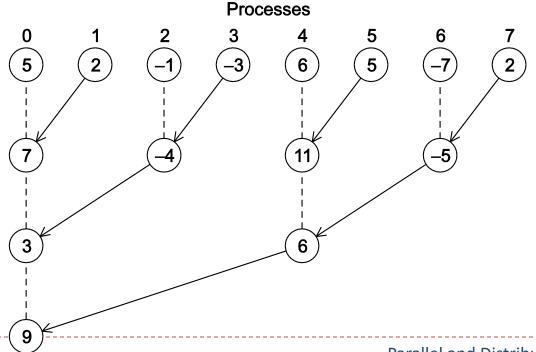
In the second phase:

- (a) Process 4 sends its newest value to process 0.
- (b) Process 0 adds the received value to its newest value.



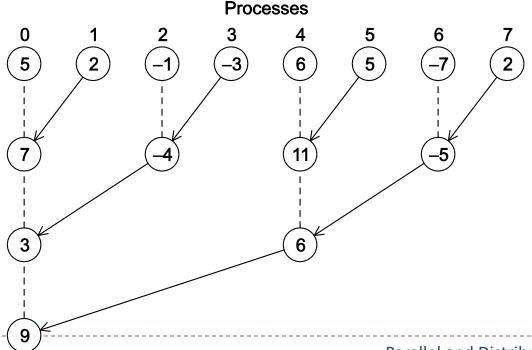
How many send/receive and add operation?

- 7 sends and 7 receives
- 7 adds



How many send/receive and add operation on master node (Rank = 0)?

- 3 sends and 3 receives
- 3 adds



Point-to-Point communication code

```
int main(void) {
                                                                      Rank = 2
      int my rank, comm sz, n = 1024, local n;
      double a = 0.0, b = 3.0, h, local a, local b;
      double local_int, total_int;
      int source:
      MPI Init(NULL, NULL);
                                                               Rank = 0 ≤
      MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
      MPI Comm size (MPI COMM WORLD, &comm sz);
10
      h = (b-a)/n; /* h is the same for all processes */
11
      local n = n/comm sz; /* So is the number of trapezoids */
12
13
14
      local_a = a + my_rank*local_n*h;
15
      local b = local a + local n*h;
      local int = Trap(local a, local b, local n, h);
16
17
      if (my_rank != 0) {
18
         MPI Send(&local int, 1, MPI DOUBLE, 0, 0,
19
20
               MPI COMM WORLD);
21
      } else {
         total int = local int;
         for (source = 1; source < comm sz; source++) {</pre>
24
            MPI_Recv(&local_int, 1, MPI_DOUBLE, source, 0,
25
                  MPI COMM WORLD, MPI STATUS IGNORE);
26
            total int += local int;
27
                                                                               uted Computation
```

Collective (tree-structure) communication

Processes

code

21

24

25

26

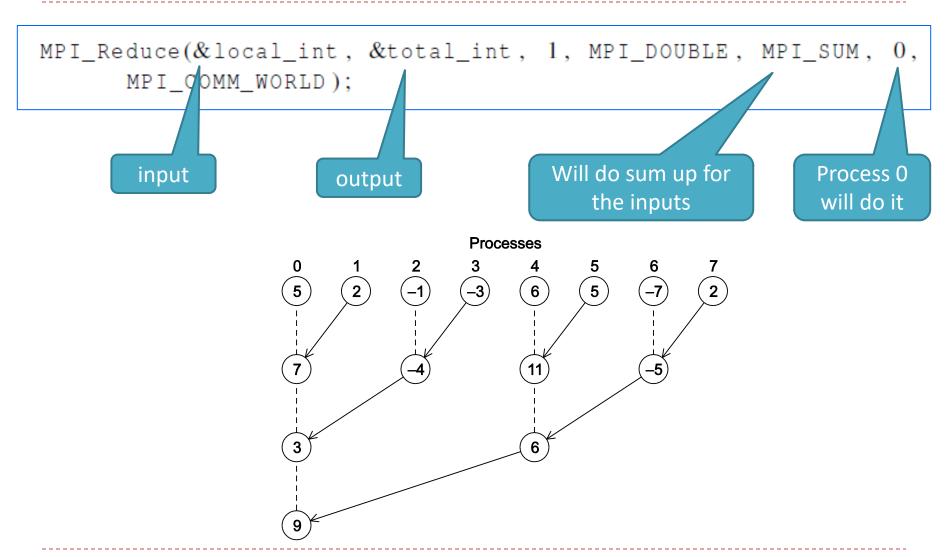
27 28

```
5
                                                                            (-3)
   int main(void) {
      int my_rank, comm_sz, n = 1024, local_n;
      double a = 0.0, b = 3.0, h, local a, local b;
      double local_int, total_int;
      int source:
      MPI_Init(NULL, NULL);
      MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
      MPI Comm size (MPI COMM WORLD, &comm sz);
10
      h = (b-a)/n; /* h is the same for all processes */
11
      local n = n/comm sz; /* So is the number of trapezoids */
12
13
      local_a = a + my_rank*local_n*h;
14
      local_b = local_a + local n*h;
15
      local int = Trap(local a, local b, local n, h);
16
17
18
      if (my_rank != 0)
         MPI Send(&local int, 1, MPI DOUBLE, 0, 0,
19
20
               MPI COMM WORLD);
                                                                 MPI Reduce(&local int,
        else {
                                                                 &total int,
        total int = local int;
                                                                 1, MPI DOUBLE,
        for (source = 1; source < comm sz; source++) {
            MPI_Recv(&local_int, 1, MPI_DOUBLE, source, 0,
                                                                 MPI SUM, 0,
                 MPI COMM WORLD, MPI STATUS IGNORE);
                                                                 MPI COMM WORLD);
            total int += local int;
```

MPI_Reduce

```
int MPI_Reduce(
     void*
                 input_data_p /* in */,
     void*
                 output_data_p /* out */,
     int
                 count
                                /* in */.
     MPI_Datatype datatype /* in */,
                            /* in */,
     qO I 9M
                 operator
                 dest_process /* in */,
     int
     MPI_Comm
                                /* in */);
                  comm
```

MPI_Reduce



Besides Sumup, it can also do:

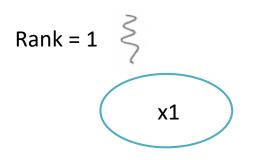
Operation Value	Meaning
MPI_MAX	Maximum
MPI_MIN	Minimum
MPI_SUM	Sum
MPI_PROD	Product
MPI_LAND	Logical and
MPI_BAND	Bitwise and
MPI_LOR	Logical or
MPI_BOR	Bitwise or
MPI_LXOR	Logical exclusive or
MPI_BXOR	Bitwise exclusive or
MPI_MAXLOC	Maximum and location of maximum
MPI_MINLOC	Minimum and location of minimum

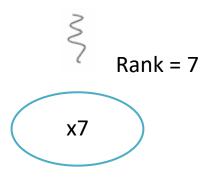
Besides Sumup, it can also do:

MPI_Reduce(&local_int, &total_int, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);

MPI_Reduce(&local_int, &total_int, 1, MPI_DOUBLE, MPI_MIN, 0, MPI_COMM_WORLD);

MPI_Reduce(&local_int, &total_int, 1, MPI_DOUBLE, MPI_PROD, 0, MPI_COMM_WORLD);

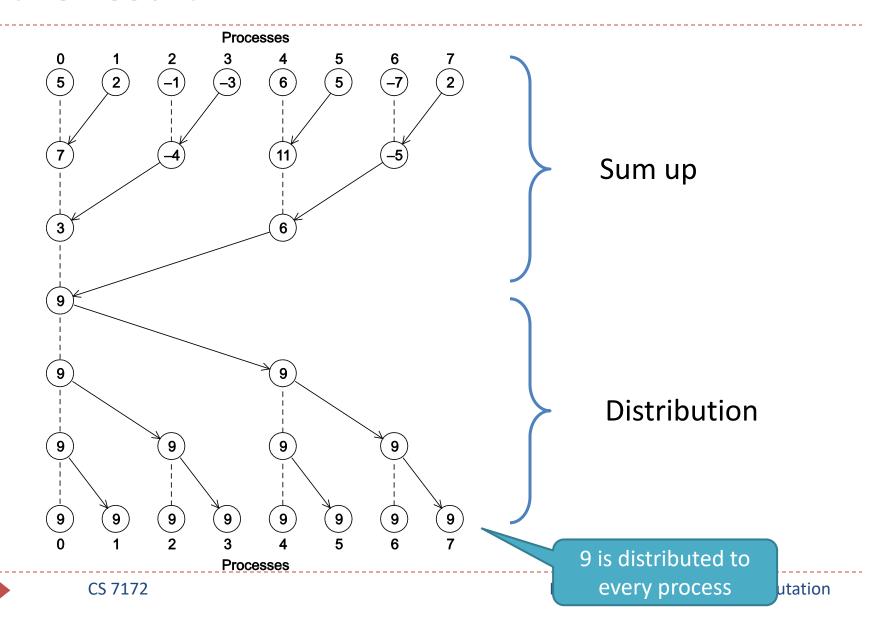




MPI_Allreduce

- Useful in a situation in which all of the processes need the final result.
 - Sum up first and then distribute to everyone

A global sum followed by distribution of the result

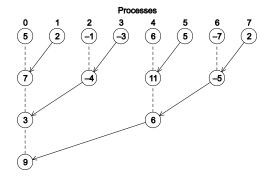


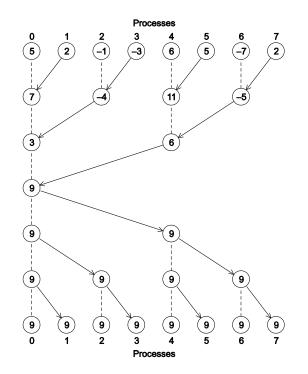
What if I only need to distribute?

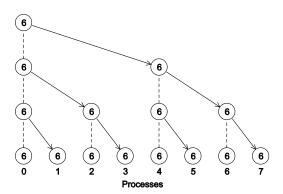
MPI_Reduce

MPI_Allreduce

?





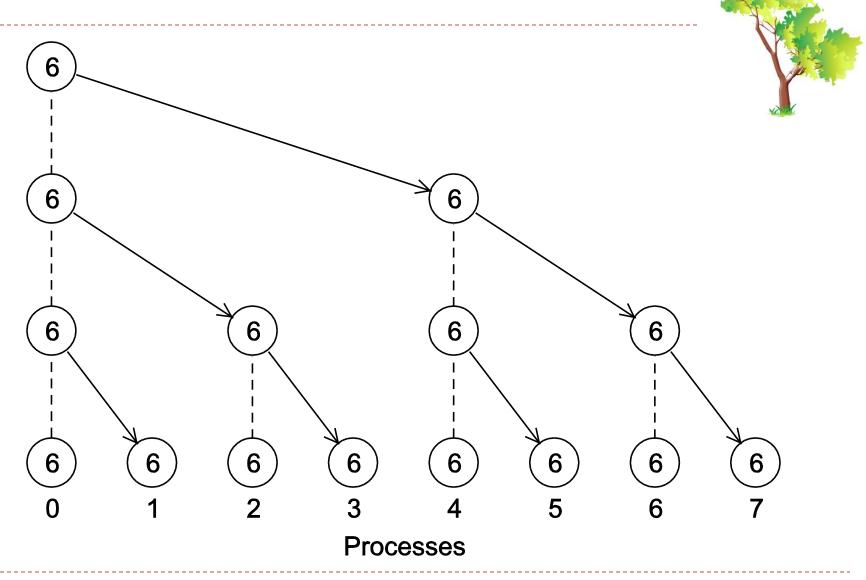


Broadcast

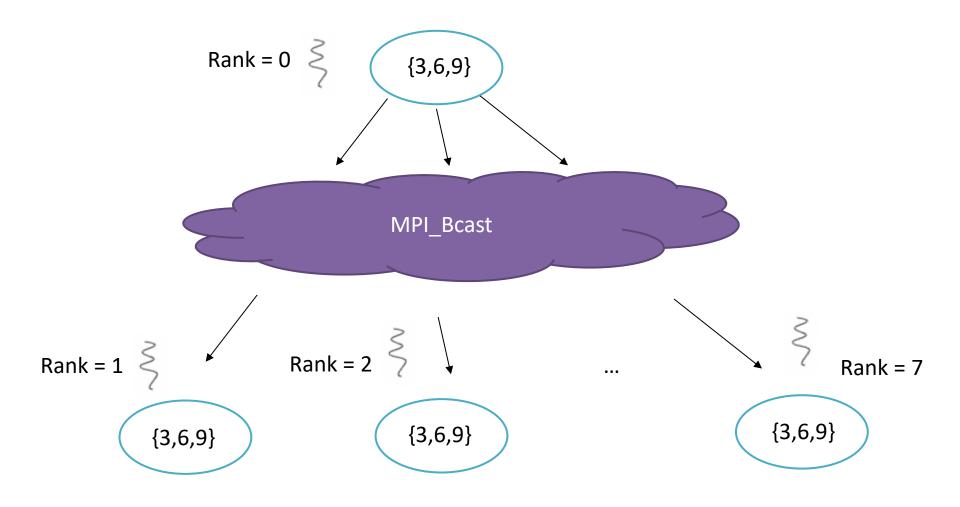
 Data belonging to a single process is sent to all of the processes in the communicator.

```
int MPI_Bcast(
        void*
                                   /* in/out
                     data p
        int
                                   /* in
                     count
        MPI_Datatype datatype
                                  /* in
        int
                                 ∕∗ in
                     source_proc
                                             */);
        MPI Comm
                                   /* in
                     comm
```

A tree-structured broadcast



Example of MPI_Bcast



Example of MPI_Bcast

```
int main(int argc, char* argv[])
{
    blog3 test;
    test.TestForMPI_Bcast(argc, argv);
}
```

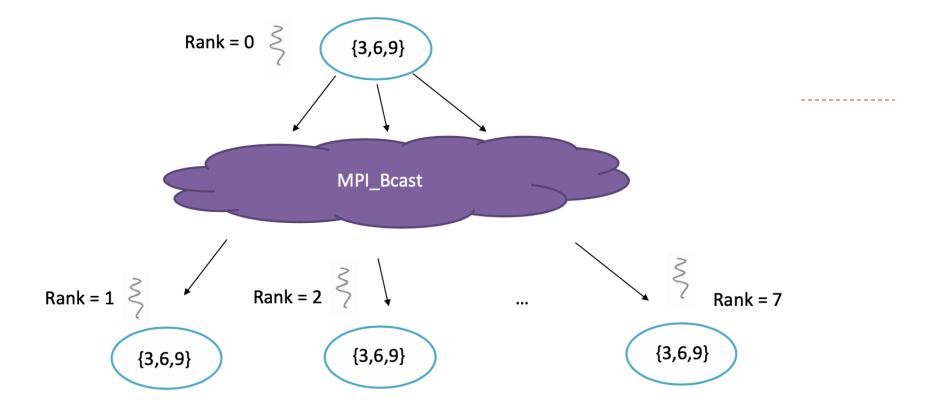
```
void blog3::TestForMPI_Bcast(int argc, char* argv[])
{
   int rankID, totalNumTasks;

   MPI_Init(&argc, &argv);
   MPI_Barrier(MPI_COMM_WORLD);
   double elapsed_time = -MPI_Wtime();

   MPI_Comm_rank(MPI_COMM_WORLD, &rankID);
   MPI_Comm_size(MPI_COMM_WORLD, &totalNumTasks);

int sendRecvBuf[3] = { 0, 0, 0 };
```

```
void blog3::TestForMPI Bcast(int argc, char* argv[])
    int rankID, totalNumTasks;
    MPI Init(&argc, &argv);
    MPI Barrier (MPI COMM WORLD);
    double elapsed time = -MPI Wtime();
    MPI Comm rank (MPI COMM WORLD, &rankID);
    MPI Comm size (MPI COMM WORLD, &totalNumTasks);
    int sendRecvBuf[3] = { 0, 0, 0 };
                                               buffer is {0,0,0}
    if (!rankID) {
        sendRecvBuf[0] = 3;
        sendRecvBuf[1] = 6;
        sendRecvBuf[2] = 9;
                                     rankID = 0 will inits
                                       buffer as {3,6,9}
    int count = 3;
    int root = 0;
   MPI Bcast (sendRecvBuf, count, MPI INT, root, MPI COMM WORLD); //MPI Bcast can be seen from all processes
    printf("my rankID = %d, sendRecvBuf = {%d, %d, %d}\n", rankID, sendRecvBuf[0], sendRecvBuf[1], sendRecvBuf[2]);
    elapsed time += MPI Wtime();
    if (!rankID) {
                                                                       Print the received data
        printf("total elapsed time = %10.6f\n", elapsed time);
    }
    MPI Finalize();
```



```
E:\CPPTest\Paiallel\vorkspace\Debug>mpiexec -n 8 Paiallel.exe

my rankID = 6, sendRecvBuf = {3, 6, 9}

my rankID = 3, sendRecvBuf = {3, 6, 9}

my rankID = 4, sendRecvBuf = {3, 6, 9}

my rankID = 7, sendRecvBuf = {3, 6, 9}

my rankID = 0, sendRecvBuf = {3, 6, 9}

total elapsed time = 0.000120

my rankID = 2, sendRecvBuf = {3, 6, 9}

my rankID = 1, sendRecvBuf = {3, 6, 9}

my rankID = 5, sendRecvBuf = {3, 6, 9}

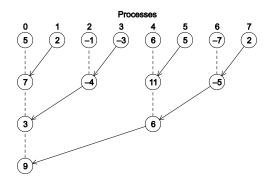
uted Computation
```

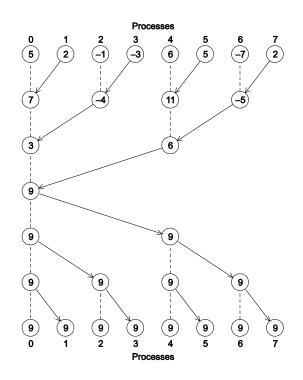
What if I only need to distribute?

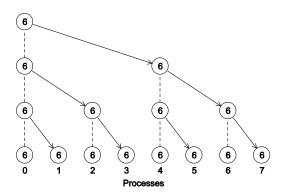
MPI_Reduce

MPI_Allreduce

MPI_Bcast







Data distributions: vector addition

я

1		-1		3		3	
0.6	+	0	+	1	=	1.6	
_2		0.2		0		-1.8	

b

1	0	+	0	-1	=	1	-1	
-1	0		-1	0		-2	0	

Partitioning options

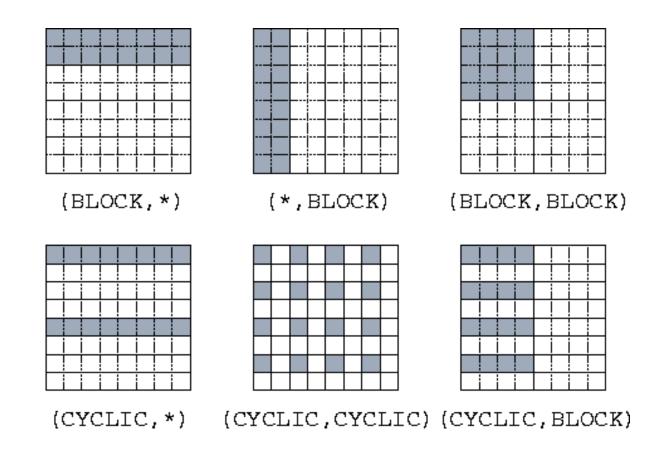
- Block partitioning
 - Assign blocks of consecutive components to each process.

- Cyclic partitioning
 - Assign components in a round robin fashion.

- Block-cyclic partitioning
 - Use a cyclic distribution of blocks of components.

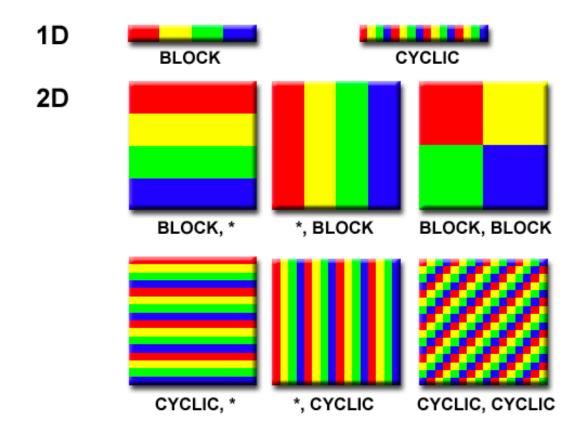
Partitioning options

Block partitioning vs Cyclic partitioning



Partitioning options

Block partitioning vs Cyclic partitioning



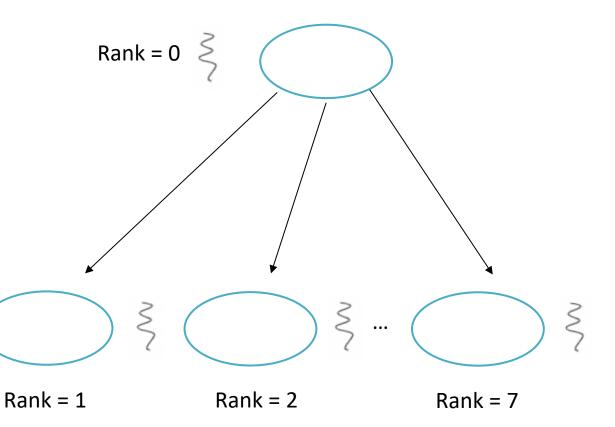
Scatter: partition & distribution

 MPI_Scatter can be used in a function that reads in an entire vector on process 0 but only sends the needed components to each of the other processes.

sendCount, each

```
process gets n data
int MPI Scatter(
     void*
                   send_buf_p
     int
                              /* in
                   send count
                               /* in
     MPI_Datatype
                   send_type
     void*
                   recv_buf_p
                               /* out */.
     int
                               /* in */.
                   recv_count
                   recv type
                               /* in */.
     MPI_Datatype
     int
                   src_proc /* in */,
                                                        scatter
                   comm
                               /* in */):
     MPI_Comm
```

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



```
int main(int argc, char* argv[])
   blog3 test;
   test.TestForMPI Scatter(argc, argv);
   return 0;
void blog3::TestForMPI Scatter(int argc, char* argv[])
   int totalNumTasks, rankID;
                                      Size = 4
   float sendBuf[SIZE][SIZE] = {
       { 1.0, 2.0, 3.0, 4.0 },
       { 5.0, 6.0, 7.0, 8.0 },
       { 9.0, 10.0, 11.0, 12.0 },
       { 13.0, 14.0, 15.0, 16.0 }
   };
   MPI Init(&argc, &argv);
   MPI Comm rank(MPI COMM WORLD, &rankID);
   MPI Comm size (MPI COMM WORLD, &totalNumTasks);
```

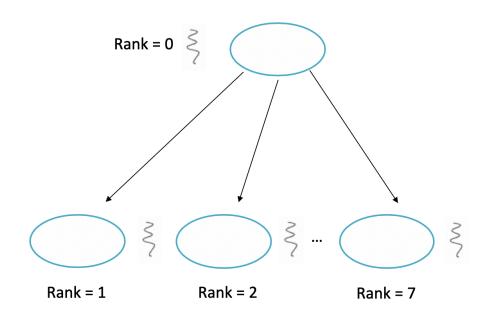
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

```
if (totalNumTasks == SIZE) {
    int source = 0;
                                    Size = 4
    int sendCount = SIZE;
   int recvCount = SIZE;
    float recvBuf[SIZE];
   //scatter data from source process to all processes in MPI COMM WORLD
   MPI Scatter (sendBuf, sendCount, MPI FLOAT,
        recvBuf, recvCount, MPI FLOAT, source, MPI COMM WORLD);
    printf("my rankID = %d, receive Results, %f %f %f %f, total = %f\n",
        rankID, recvBuf[0], recvBuf[1], recvBuf[2], recvBuf[3],
        recvBuf[0] + recvBuf[1] + recvBuf[2] + recvBuf[3]);
else if (totalNumTasks == 8) {
                                      sendCount = 2, each
    int source = 0;
    int sendCount = 2;
                                       process gets 2 data
   int recvCount = 2;
    float recvBuf[2];
   MPI Scatter (sendBuf, sendCount, MPI FLOAT,
        recvBuf, recvCount, MPI FLOAT, source, MPI COMM WORLD);
    printf("my rankID = %d, receive result: %f %f, total = %f\n",
        rankID, recvBuf[0], recvBuf[1], recvBuf[0] + recvBuf[1]);
else {
    printf("Please specify -n %d or -n %d\n", SIZE, 2 * SIZE);
MPI Finalize();
```

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

After each process receives the data, sum them up

I and Distributed Computation



```
      1
      2
      3
      4

      5
      6
      7
      8

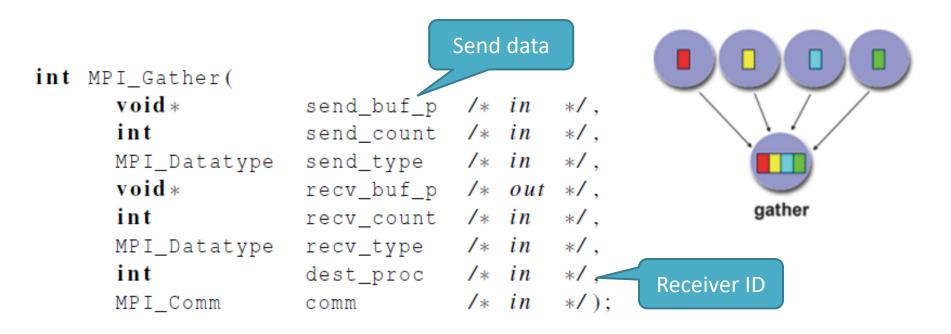
      9
      10
      11
      12

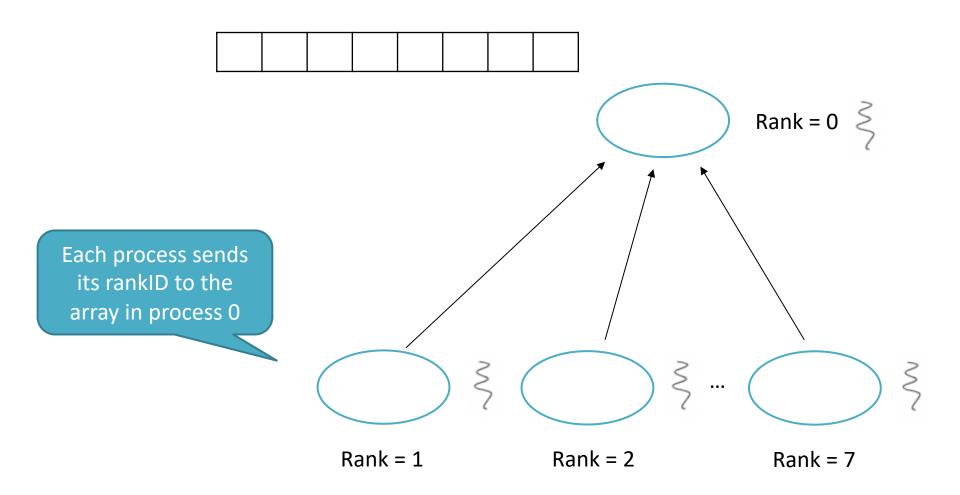
      13
      14
      15
      16
```

```
E:\CPPTest\Paiallel\vorkspace\Debug>mpiexec -n 8 Paiallel.exe
my rankID = 3, receive result: 7.000000 8.000000, total = 15.000000
my rankID = 5, receive result: 11.000000 12.000000, total = 23.000000
my rankID = 2, receive result: 5.000000 6.000000, total = 11.000000
my rankID = 0, receive result: 1.000000 2.000000, total = 3.000000
my rankID = 7, receive result: 15.000000 16.000000, total = 31.000000
my rankID = 4, receive result: 9.000000 10.000000, total = 19.000000
my rankID = 1, receive result: 3.000000 4.000000, total = 7.000000
my rankID = 6, receive result: 13.000000 14.000000, total = 27.000000
```

Gather

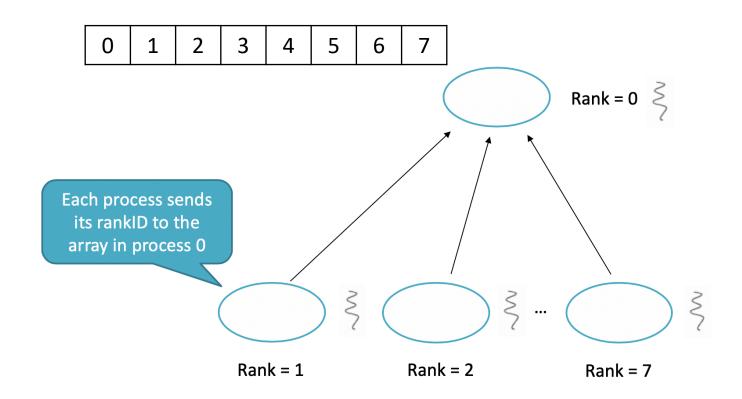
 Collect all of the components of the vector onto process 0, and then process 0 can process all of the components.





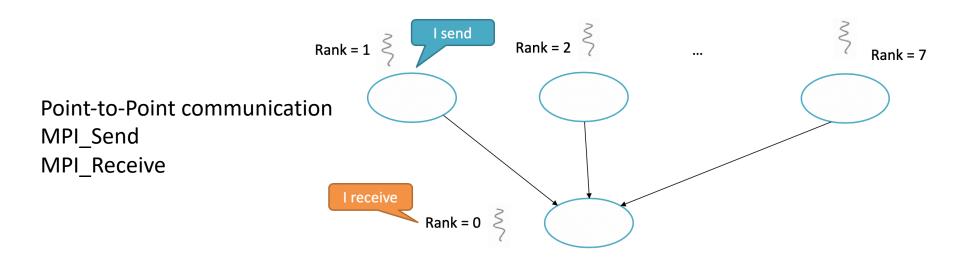
```
int main(int argc, char* argv[])
    blog3 test;
    test.TestForMPI Gather(argc, argv);
    return 0;
void blog3::TestForMPI Gather(int argc, char* argv[])
    int rankID, totalNumTasks;
    MPI Init(&argc, &argv);
    MPI Barrier (MPI COMM WORLD);
    double elapsed time = -MPI Wtime();
    MPI Comm rank (MPI COMM WORLD, &rankID);
    MPI Comm size (MPI COMM WORLD, &totalNumTasks);
```

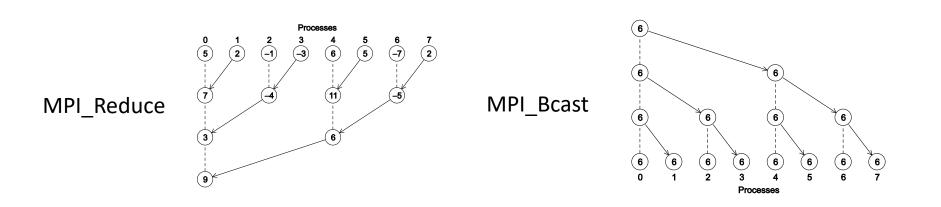
```
int* gatherBuf = (int *)malloc(sizeof(int) * totalNumTasks);
if (gatherBuf == NULL) {
    printf("malloc error!");
    exit(-1);
                                    The sent content
    MPI Finalize();
                                      is the rankID
int sendBuf = rankID; //for eac
                                   ocess, its rankID will be sent out
                                                                               Receiver is
                                                                                process 0
int sendCount = 1;
                                  Data size is 1
int recvCount = 1;
int root = 0;
MPI Gather (&sendBuf, sendCount, MPI INT, gatherBuf, recvCount, MPI INT, root, MPI COMM WORLD);
elapsed time += MPI Wtime();
if (!rankID) {
    int i;
    for (i = 0; i < totalNumTasks; i++) {</pre>
        printf("gatherBuf[%d] = %d, ", i, gatherBuf[i]);
    putchar('\n');
    printf("total elapsed time = %10.6f\n", elapsed time);
MPI Finalize();
```



```
E:\CPPTest\Paiallel\orkspace\Debug>mpiexec -n 8 Paiallel.exe
gatherBuf[0] = 0, gatherBuf[1] = 1, gatherBuf[2] = 2, gatherBuf[3] = 3, gatherBuf[4] = 4, gatherB
uf[5] = 5, gatherBuf[6] = 6, gatherBuf[7] = 7,
total elapsed time = 0.000850
```

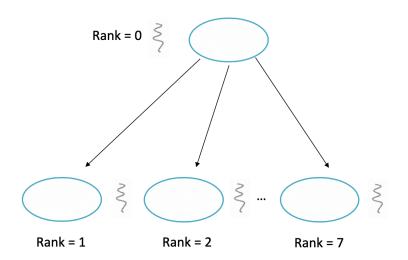
Conclusion

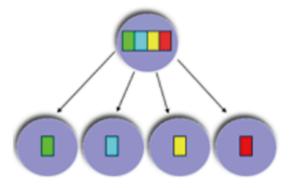




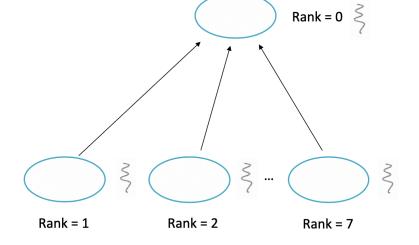
Conclusion

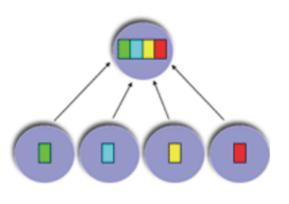
MPI_Scatter



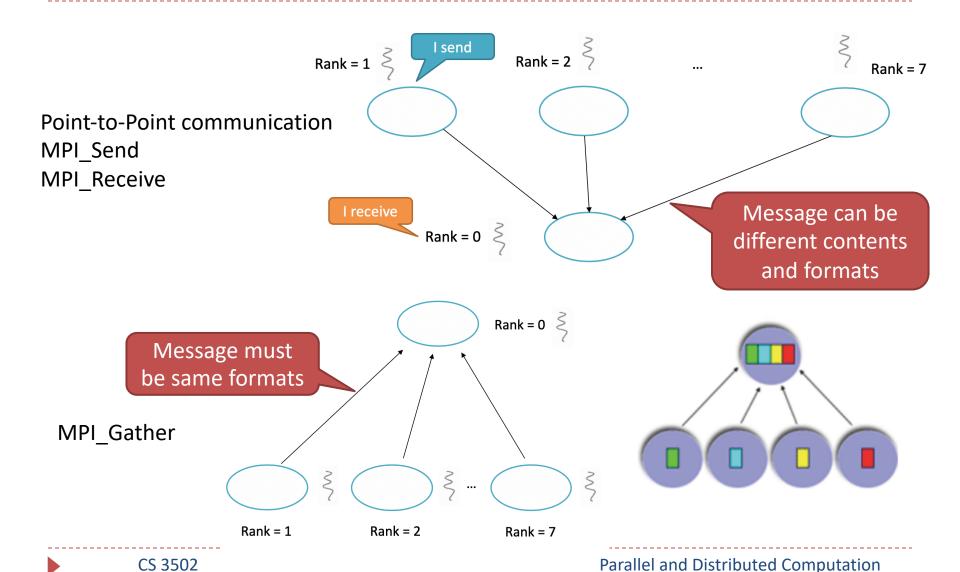


MPI_Gather





What is the Difference between Point-to-Point communication and MPI_Gather?



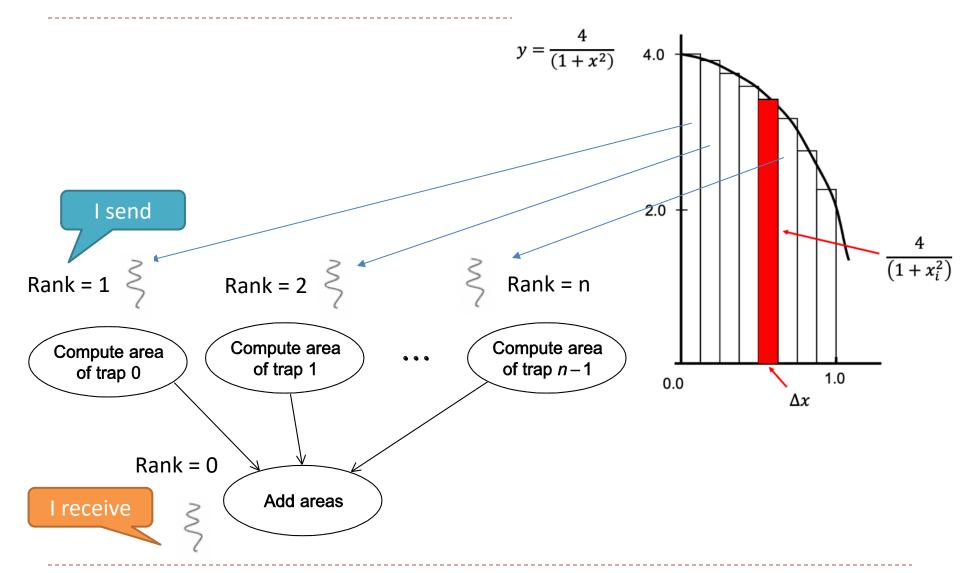
$\int_0^1 \frac{4}{(1+x^2)} dx = \pi$

Hint of Project 1

```
#define NUMSTEPS 1000000
int main() {
        int i;
        double x, pi, sum = 0.0;
        struct timespec start, end;
        clock gettime(CLOCK MONOTONIC, &start);
                                                                     2.0
        double step <1.0/(double) NUMSTEPS;
        x = 0.5 * step;
        for (i=0;i<= NUMSTEPS; i++) {</pre>
                x+=step;
                sum +=(4.0/(1.0+x*x);)
                                                                                          1.0
                                                                       0.0
        pi = step * sum; size of rectangle
                                                                                      \Delta x
        clock gettime (CLOCK MONOTONIC, &end);
        u int64 t diff = 100000000L * (end.tv sec - start.tv sec) + end.tv nsec -
start.tv nsec;
        printf("PI is %.20f\n",pi);
        printf("elapsed time = %llu nanoseconds\n", (long long unsigned int) diff);
        return 0;
```

Hint of Project 1

$$\int_0^1 \frac{4}{(1+x^2)} dx = \pi$$



Hint of Project 1

$$\int_0^1 \frac{4}{(1+x^2)} dx = \pi$$

