

# Communicating NumPy arrays

Xin Li

# Overview

- Numpy and buffer-like object
- Blocking communication
- Nonblocking communication
- Collective communication

Numpy and buffer-like object

# NumPy

- Powerful N-dimensional arrays
- Optimized for performance
- Easy to use
- Open source

# NumPy

```
>>> import numpy as np
```

```
>>> a = np.arange(16.).reshape(4,4)
```

```
>>> a
```

```
array([[ 0.,  1.,  2.,  3.],  
       [ 4.,  5.,  6.,  7.],  
       [ 8.,  9., 10., 11.],  
       [12., 13., 14., 15.]])
```

```
>>> b = a.T
```

```
>>> b
```

```
array([[ 0.,  4.,  8., 12.],  
       [ 1.,  5.,  9., 13.],  
       [ 2.,  6., 10., 14.],  
       [ 3.,  7., 11., 15.]])
```

# NumPy

```
>>> import numpy as np
```

```
>>> a = np.arange(16.).reshape(4,4)
```

```
>>> a
```

```
array([[ 0.,  1.,  2.,  3.],  
       [ 4.,  5.,  6.,  7.],  
       [ 8.,  9., 10., 11.],  
       [12., 13., 14., 15.]])
```

```
>>> c = a[0:2, 0:2]
```

```
>>> c
```

```
array([[0., 1.],  
       [4., 5.]])
```

# NumPy

```
>>> import numpy as np
```

```
>>> a = np.arange(16.).reshape(4,4)
```

```
>>> a
```

```
array([[ 0.,  1.,  2.,  3.],  
       [ 4.,  5.,  6.,  7.],  
       [ 8.,  9., 10., 11.],  
       [12., 13., 14., 15.]])
```

```
>>> d = a[:,1].reshape(2,2)
```

```
>>> d
```

```
array([[ 1.,  5.],  
       [ 9., 13.]])
```

# NumPy

```
>>> import numpy as np
```

```
>>> a = np.arange(16.).reshape(4,4)
```

```
>>> a
```

```
array([[ 0.,  1.,  2.,  3.],  
       [ 4.,  5.,  6.,  7.],  
       [ 8.,  9., 10., 11.],  
       [12., 13., 14., 15.]])
```

```
>>> np.matmul(a, a.T)
```

```
array([[14., 38., 62., 86.],  
       [38., 126., 214., 302.],  
       [62., 214., 366., 518.],  
       [86., 302., 518., 734.]])
```



# NumPy with mpi4py

- Numpy arrays can be communicated by the all-lowercase methods like send, recv, bcast, etc.
- For best efficiency Numpy arrays can be communicated as **buffer-like objects**, using method with a leading uppercase letter.
  - Send, Recv
  - Isend, Irecv
  - Bcast, Reduce
  - Scatterv, Gatherv (more useful than Scatter/Gather)

# Numpy array as buffer-like objects

- communication is fast
  - close to the speed of MPI communication in C
- less flexible
  - memory of the receiving buffer needs to be allocated
  - size of the sending buffer should not exceed that of the receiving buffer
  - mpi4py expects the buffer-like objects to have contiguous memory

# Contiguous or non-contiguous?

```
>>> import numpy as np
```

```
>>> a = np.arange(16.).reshape(4,4)
```

```
>>> a
```

```
array([[ 0.,  1.,  2.,  3.],  
       [ 4.,  5.,  6.,  7.],  
       [ 8.,  9., 10., 11.],  
       [12., 13., 14., 15.]])
```

```
>>> a.flags
```

```
  C_CONTIGUOUS : True
```

```
  F_CONTIGUOUS : False
```

```
  OWNDATA : False
```

```
...
```

# Contiguous or non-contiguous?

```
>>> b = a.T
>>> b.flags
  C_CONTIGUOUS : False
  F_CONTIGUOUS : True
  OWNDATA : False
  ...

>>> b[0,0] = 99.0
>>> b
array([[99.,  4.,  8., 12.],
       [ 1.,  5.,  9., 13.],
       [ 2.,  6., 10., 14.],
       [ 3.,  7., 11., 15.]])
```

# Contiguous or non-contiguous?

```
>>> b = a.T
>>> b.flags
  C_CONTIGUOUS : False
  F_CONTIGUOUS : True
  OWNDATA : False
  ...

>>> b[0,0] = 99.0
>>> b
array([[99.,  4.,  8., 12.],
       [ 1.,  5.,  9., 13.],
       [ 2.,  6., 10., 14.],
       [ 3.,  7., 11., 15.]])

>>> a
array([[99.,  1.,  2.,  3.],
       [ 4.,  5.,  6.,  7.],
       [ 8.,  9., 10., 11.],
       [12., 13., 14., 15.]])
```

# Contiguous or non-contiguous?

```
>>> a
array([[99.,  1.,  2.,  3.],
       [ 4.,  5.,  6.,  7.],
       [ 8.,  9., 10., 11.],
       [12., 13., 14., 15.]])
```

```
>>> c = a[0:2, 0:2]
```

```
>>> c
array([[99.,  1.],
       [ 4.,  5.]])
```

```
>>> c.flags
C_CONTIGUOUS : False
F_CONTIGUOUS : False
OWNDATA : False
```

# Contiguous or non-contiguous?

```
>>> a
array([[99.,  1.,  2.,  3.],
       [ 4.,  5.,  6.,  7.],
       [ 8.,  9., 10., 11.],
       [12., 13., 14., 15.]])
```



```
>>> c = a[0:2, 0:2]
>>> c
array([[99.,  1.],
       [ 4.,  5.]])
```



```
>>> c.flags
C_CONTIGUOUS : False
F_CONTIGUOUS : False
OWNDATA : False
```

# Contiguous or non-contiguous?

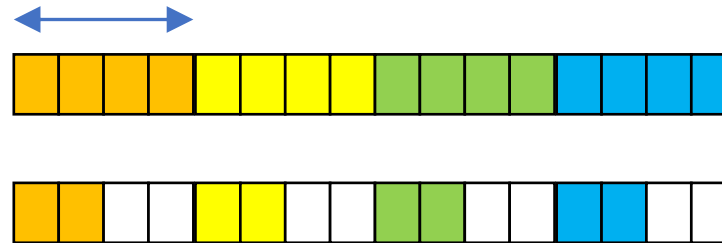
```
>>> d = c.copy()
>>> d.flags
  C_CONTIGUOUS : True
  F_CONTIGUOUS : False
  OWNDATA : True
  ...
```



# Contiguous or non-contiguous?

```
>>> d = c.copy()
>>> d.flags
  C_CONTIGUOUS : True
  F_CONTIGUOUS : False
  OWNDATA : True
  ...
```

```
>>> a.strides
(32, 8)
>>> c.strides
(32, 8)
>>> d.strides
(16, 8)
```



# Use NumPy array as buffer-like object

- The array itself, or a list or tuple with
  - 2 or 3 elements
  - 4 elements for the vector variants (Scatterv, Gatherv)
- data
- [data, MPI.DOUBLE]
- [data, n, MPI.DOUBLE]
- [data, count, displ, MPI.DOUBLE]

Blocking send/recv

# Send / Recv

- Syntax
  - `comm.Send(obj, dest=dest, tag=tag)`
  - `comm.Recv(obj, source=source, tag=rag)`
- Note
  - obj needs to be created prior to the communication
  - size of obj needs to be known before hand

# Send / Recv

- example

```
from mpi4py import MPI
import numpy as np
```

```
comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()
```

```
if rank == 0:
    data = np.arange(4.)
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(4)
    comm.Recv(data, source=0, tag=rank)
    print('Process {} received data:'.format(rank), data)
```

# Send / Recv

- output

Process 0 sent data: [0. 1. 2. 3.]

Process 0 sent data: [0. 1. 2. 3.]

Process 0 sent data: [0. 1. 2. 3.]

Process 1 received data: [0. 1. 2. 3.]

Process 2 received data: [0. 1. 2. 3.]

Process 3 received data: [0. 1. 2. 3.]

# Send / Recv

- question
  - If the size of the numpy array is only known on the master process (rank 0), how do we send it to other processes (rank > 0)?

# Send / Recv

- question
  - If the size of the numpy array is only known on the master process (rank 0), how do we send it to other processes (rank > 0)?
- solution
  - need to send the size of the array first
- exercise
  - rewrite the above example assuming that the size of array is not known on the non-master nodes (rank > 0).
  - for sending an integer you may use the all-lowercase send



# Send / Recv

- exercise
  - what will happen if you try to send an array with non-contiguous memory?
  - hint: this is how to create a simple array with non-contiguous memory  
`data = np.arange(12.)[:, :2]`

# Send / Recv

- exercise
  - what will happen if the receiving buffer is ***larger*** than the sent array?

```
if rank == 0:
    data = np.arange(4.)
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else
    data = np.zeros(6)
    comm.Recv(data, source=0, tag=rank)
    print('Process {} has data:'.format(rank), data)
```

# Send / Recv

- exercise
  - what will happen if the receiving buffer is *smaller* than the sent array?

```
if rank == 0:
    data = np.arange(4.)
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else
    data = np.zeros(3)
    comm.Recv(data, source=0, tag=rank)
    print('Process {} has data:'.format(rank), data)
```

# Send / Recv

- best practice: check status

```
if rank == 0:
    data = np.arange(4.)
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(3)
    status = MPI.Status()
    comm.Recv(data, source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('Process {} received data:'.format(rank), data)
```

# Send / Recv with buffer size

- use [data, n, MPI.DOUBLE] to specify the buffer

```
if rank == 0:
    data = np.arange(4.)
    for i in range(1, size):
        comm.Send([data, 2, MPI.DOUBLE], dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data[:2])
else:
    data = np.zeros(4)
    status = MPI.Status()
    comm.Recv([data, 2, MPI.DOUBLE], source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('Process {} received data:'.format(rank), data[:2])
```

# Send / Recv with buffer size

- use [data, n, MPI.DOUBLE] to specify the buffer

```
if rank == 0:
    data = np.arange(4.)
    for i in range(1, size):
        comm.Send([data, 2, MPI.DOUBLE], dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data[:2])
else:
    data = np.zeros(4)
    status = MPI.Status()
    comm.Recv([data, 2, MPI.DOUBLE], source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('Process {} received data:'.format(rank), data[:2])
```

- NOTE: The size of buffer should never be larger than the size of the Numpy array.

# Send / Recv with buffer size

- use array slicing

```
if rank == 0:
    data = np.arange(10.)
    for i in range(1, size):
        comm.Send(data[2:6], dest=i, tag=i)
    print('data on process {}:'.format(rank), data)
else:
    data = np.zeros(10)
    status = MPI.Status()
    comm.Recv(data[2:6], source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('data on process {}:'.format(rank), data)
```

# Send / Recv with 2D array

- exercise

- Send / Recv 2D array

```
if rank == 0:
    data = np.arange(16.).reshape(4,4)
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(16).reshape(4,4)
    status = MPI.Status()
    comm.Recv(data, source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('Process {} received data:'.format(rank), data)
```

- Will it work if the receiving buffer is a 1D array?



# Send / Recv with 2D array

- exercise

- What will happen if we send a 2D array with .T (transpose)?

```
if rank == 0:
    data = np.arange(16.).reshape(4,4).T
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(16).reshape(4,4)
    status = MPI.Status()
    comm.Recv(data, source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('Process {} received data:'.format(rank), data)
```

# Send / Recv with 2D array

- exercise
  - What will happen if we send a 2D array with `.T.copy()` (copy of transpose)?

```
if rank == 0:
    data = np.arange(16.).reshape(4,4).T.copy()
    for i in range(1, size):
        comm.Send(data, dest=i, tag=i)
        print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(16).reshape(4,4)
    status = MPI.Status()
    comm.Recv(data, source=0, tag=rank, status=status)
    if status.error != 0:
        comm.Abort(status.error)
    print('Process {} received data:'.format(rank), data)
```

Non-blocking send/rcv

# Isend / Irecv

- Syntax

- `comm.Isend(obj, dest=dest, tag=tag)`
- `comm.Irecv(obj, source=source, tag=rag)`

- Note

- obj needs to be created prior to the communication
- size of obj needs to be known before hand
- A Request object is returned by Isend / Irecv
- Use Wait method of the Request object
- No Status involved

# Isend / Irecv

- example

```
from mpi4py import MPI
import numpy as np
```

```
comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()
```

```
if rank == 0:
    data = np.arange(4.)
    reqs = []
    for i in range(1, size):
        reqs.append(comm.Isend(data, dest=i, tag=i))
    for req in reqs:
        req.wait()
    print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(4)
    req = comm.Irecv(data, source=0, tag=rank)
    req.wait()
    print('Process {} received data:'.format(rank), data)
```

# Isend / Irecv

- example

```
from mpi4py import MPI
import numpy as np
```

```
comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()
```

```
if rank == 0:
    data = np.arange(4.)
    reqs = []
    for i in range(1, size):
        reqs.append(comm.Isend(data, dest=i, tag=i))
    for req in reqs:
        req.wait()
    print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(4)
    req = comm.Irecv(data, source=0, tag=rank)
    req.wait()
    print('Process {} received data:'.format(rank), data)
```

# Isend / Irecv

- example

```
from mpi4py import MPI
import numpy as np
```

```
comm = MPI.COMM_WORLD
rank = comm.Get_rank()
size = comm.Get_size()
```

```
if rank == 0:
    data = np.arange(4.)
    reqs = [comm.Isend(data, dest=i, tag=i) for i in range(1, size)]
    for req in reqs:
        req.wait()
        print('Process {} sent data:'.format(rank), data)
else:
    data = np.zeros(4)
    req = comm.Irecv(data, source=0, tag=rank)
    req.wait()
    print('Process {} received data:'.format(rank), data)
```

# Isend / Irecv

- exercise:
  - what if the size of the receiving buffer is larger than the sent array?
  - what if the size of the receiving buffer is smaller than the sent array?
- send a 2D array
  - without transpose
  - with .T
  - with .T.copy()



Collectives

# Bcast

- Syntax
  - `comm.Bcast(obj, root=root)`
- Note
  - obj needs to be created prior to the communication
  - size of obj needs to be known before hand
  - root can be non-zero (source of communication)

# Bcast

```
from mpi4py import MPI
import numpy as np

comm = MPI.COMM_WORLD
rank = comm.Get_rank()

if rank == 0:
    data = np.arange(4.0)
else:
    data = np.zeros(4)

comm.Bcast(data, root=0)

print('Process {} has data:'.format(rank), data)
```

# Bcast

- exercise
  - what will happen if the size of the receiving buffer is ***larger***?
  - what will happen if the size of the receiving buffer is ***smaller***?

# Scatterv

- vector variant of Scatter
- Syntax
  - `comm.Scatterv([sendbuf, count, displ, MPI.DOUBLE], recvbuf, root=root)`
- Note
  - `[sendbuf, count, displ, MPI.DOUBLE]` defines the sending buffer
  - `recvbuf` needs to be created prior to the communication
  - size of `recvbuf` needs to be known before hand
  - `root` can be non-zero (source of Scatterv)

# Scatterv

- determine count and displ

```
ave, res = divmod(sendbuf.size, nprocs)

# count: the size of each sub-task
count = np.array([ave + 1 if p < res else ave for p in range(nprocs)])

# displacement: the starting index of each sub-task
displ = np.array([sum(count[:p]) for p in range(nprocs)])
```

# Scatterv

- determine count and displ

```
ave, res = divmod(sendbuf.size, nprocs)

# count: the size of each sub-task
count = np.array([ave + 1 if p < res else ave for p in range(nprocs)])

# displacement: the starting index of each sub-task
displ = np.array([sum(count[:p]) for p in range(nprocs)])
```

- `sendbuf.size = 15; nprocs = 4`
  - `ave = 3; res = 3`
  - `count = np.array([4, 4, 4, 3])`
  - `displ = np.array([0, 4, 8, 12])`

# Scatterv

- determine count and displ

```
ave, res = divmod(sendbuf.size, nprocs)
```

```
# count: the size of each sub-task
```

```
count = np.array([ave + 1 if p < res else ave for p in range(nprocs)])
```

```
# displacement: the starting index of each sub-task
```

```
displ = np.array([sum(count[:p]) for p in range(nprocs)])
```

- sendbuf.size = 15; nprocs = 4

- ave = 3; res = 3
- count = np.array([4, 4, 4, 3])
- displ = np.array([0, 4, 8, 12])

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

rank 0

rank 1

rank 2

rank 3



# Scatterv

```
from mpi4py import MPI  
import numpy as np
```

```
comm = MPI.COMM_WORLD  
rank = comm.Get_rank()  
nprocs = comm.Get_size()
```

# Scatterv

```
if rank == 0:
    sendbuf = np.arange(15.0)

    # count: the size of each sub-task
    ave, res = divmod(sendbuf.size, nprocs)
    count = [ave + 1 if p < res else ave for p in range(nprocs)]
    count = np.array(count)

    # displacement: the starting index of each sub-task
    displ = [sum(count[:p]) for p in range(nprocs)]
    displ = np.array(displ)

else:
    sendbuf = None

    # initialize count on worker processes
    count = np.zeros(nprocs, dtype=np.int)
    displ = None
```

# Scatterv

```
# broadcast count
```

```
comm.Bcast(count, root=0)
```

```
# initialize recvbuf on all processes
```

```
recvbuf = np.zeros(count[rank])
```

```
comm.Scatterv([sendbuf, count, displ, MPI.DOUBLE], recvbuf, root=0)
```

```
print('After Scatterv, process {} has data:'.format(rank), recvbuf)
```

# Scatterv

- output

After Scatterv, process 0 has data: [0. 1. 2. 3.]

After Scatterv, process 2 has data: [ 8. 9. 10. 11.]

After Scatterv, process 1 has data: [4. 5. 6. 7.]

After Scatterv, process 3 has data: [12. 13. 14.]

# Gatherv

- vector variant of Gather
- Syntax
  - `comm.Gather(sendbuf, [recvbuf, count, displ, MPI.DOUBLE], root=root)`
- Note
  - `[recvbuf, count, displ, MPI.DOUBLE]` defines the receiving buffer
  - `recvbuf` needs to be created prior to the communication
  - size of `recvbuf` needs to be known before hand
  - `root` can be non-zero (target of Gatherv)

# Gatherv

- continue with the Scatterv code

```
sendbuf2 = recvbuf  
recvbuf2 = np.zeros(sum(count))
```

```
comm.Gatherv(sendbuf2, [recvbuf2, count, displ, MPI.DOUBLE], root=0)
```

```
if comm.Get_rank() == 0:  
    print('After Gatherv, process 0 has data:', recvbuf2)
```

# Gatherv

- output

After Gatherv, process 0 has data:

```
[ 0.  1.  2.  3.  4.  5.  6.  7.  8.  9. 10. 11. 12. 13. 14.]
```

# Reduce

- Syntax
  - `comm.Reduce(sendbuf, recvbuf, op=op, root=root)`
- Note
  - default op is `MPI.SUM`
  - root can be non-zero (target of Reduce)



# Reduce

- continue with the Scatterv code

```
partial_sum = np.zeros(1)
partial_sum[0] = sum(recvbuf)
print('Partial sum on process {} is:'.format(rank), partial_sum[0])

total_sum = np.zeros(1)
comm.Reduce(partial_sum, total_sum, op=MPI.SUM, root=0)

if comm.Get_rank() == 0:
    print('After Reduce, total sum on process 0 is:', total_sum[0])
```

# Reduce

- output

```
Partial sum on process 3 is: 39.0
```

```
Partial sum on process 1 is: 22.0
```

```
Partial sum on process 2 is: 38.0
```

```
Partial sum on process 0 is: 6.0
```

```
After Reduce, total sum on process 0 is: 105.0
```

# Reduce

- exercise
  - Use Reduce to compute the sum of numpy arrays

# Reduce

```
from mpi4py import MPI
import numpy as np

comm = MPI.COMM_WORLD
rank = comm.Get_rank()
nprocs = comm.Get_size()

partial = np.arange(10.) * rank
print('Before Reduce, partial on process {} is:'.format(rank), partial)

total = np.zeros(10)
comm.Reduce(partial, total, op=MPI.SUM, root=0)
if comm.Get_rank() == 0:
    print('After Reduce, total on process 0 is:', total)
```

# Reduce

```
Before Reduce, partial on process 0 is: [0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
Before Reduce, partial on process 2 is: [ 0.  2.  4.  6.  8. 10. 12. 14. 16. 18.]
Before Reduce, partial on process 1 is: [0. 1. 2. 3. 4. 5. 6. 7. 8. 9.]
Before Reduce, partial on process 3 is: [ 0.  3.  6.  9. 12. 15. 18. 21. 24. 27.]
After Reduce, total on process 0 is: [ 0.  6. 12. 18. 24. 30. 36. 42. 48. 54.]
```

# Summary

# Numpy array as buffer-like object in mpi4py

- fast communication
- less flexible code (need to deal with memory)
- blocking and nonblocking communication
  - Send, Recv, Isend, Irecv
- collective communication
  - Bcast, Scatterv, Gatherv, Reduce