

# MPI for Python

<http://mpi4py.scipy.org>

Lisandro Dalcin

[dalcinl@gmail.com](mailto:dalcinl@gmail.com)

Centro Internacional de Métodos Computacionales en Ingeniería  
Consejo Nacional de Investigaciones Científicas y Técnicas  
Santa Fe, Argentina

January, 2011

Python for parallel scientific computing  
PASI, Valparaíso, Chile

# Outline

Overview

Communicators

Point to Point Communication

Collective Operations

Compute Pi

Mandelbrot Set

Dynamic Process Management

## Overview

Communicators

Point to Point Communication

Collective Operations

Compute Pi

Mandelbrot Set

Dynamic Process Management

# What is `mpi4py`?

- ▶ Full-featured Python bindings for **MPI**.
- ▶ API based on the standard MPI-2 C++ bindings.
- ▶ Almost all MPI calls are supported.
  - ▶ targeted to MPI-2 implementations.
  - ▶ also works with MPI-1 implementations.

# Implementation

## Implemented with **Cython**

- ▶ Code base far easier to write, maintain, and extend.
- ▶ Faster than other solutions (mixed Python and C codes).
- ▶ A *pythonic* API that runs at C speed !

# Features – MPI-1

- ▶ Process groups and communication domains.
  - ▶ intracommunicators
  - ▶ intercommunicators
- ▶ Point to point communication.
  - ▶ blocking (send/rcv)
  - ▶ nonblocking (isend/irecv + test/wait)
- ▶ Collective operations.
  - ▶ Synchronization (barrier)
  - ▶ Communication (broadcast, scatter/gather)
  - ▶ Global reductions (reduce, scan)

## Features – MPI-2

- ▶ Dynamic process management (spawn, connect/accept).
- ▶ Parallel I/O (read/write).
- ▶ One sided operations, a.k.a. RMA (put/get/accumulate).
- ▶ Extended collective operations.

# Features – Python

- ▶ Communication of Python objects.
  - ▶ high level and very convenient, based in `pickle` serialization
  - ▶ can be slow for large data (CPU and memory consuming)

```
<object> → pickle.dump() → send()
```



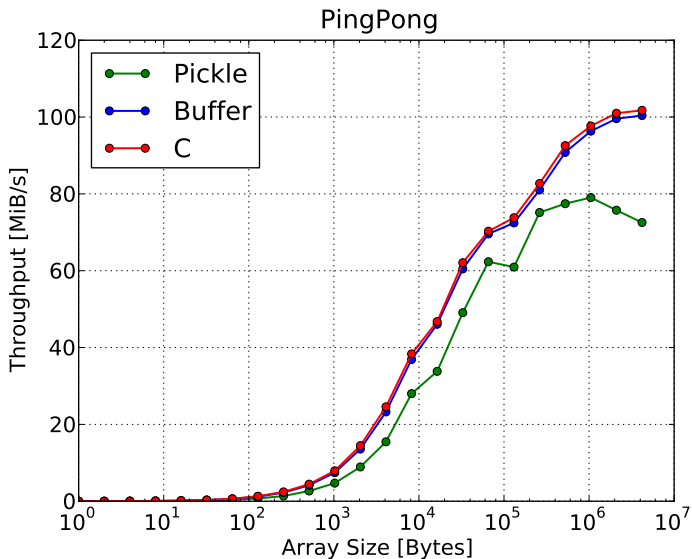
```
<object> ← pickle.load() ← recv()
```

- ▶ Communication of array data (e.g. **NumPy** arrays).
  - ▶ lower level, slightly more verbose
  - ▶ very fast, almost C speed (for messages above 5-10 KB)

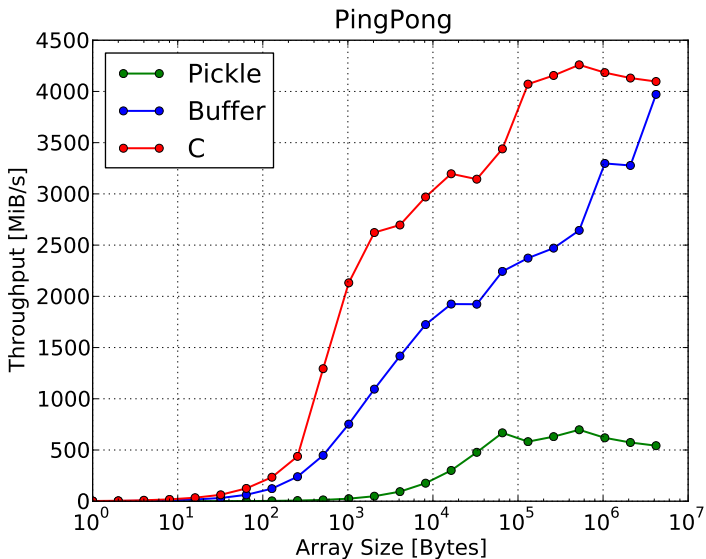
```
message = [<object>, (count, displ), datatype]
```



# Point to Point Throughput – Gigabit Ethernet



# Point to Point Throughput – Shared Memory



## Features – IPython

Integration with **IPython** enables MPI to be used *interactively*.

- ▶ Start engines with MPI enabled

```
$ ipcluster mpiexec -n 16 --mpi=mpi4py
```

- ▶ Connect to the engines

```
$ ipython
```

```
In [1]: from IPython.kernel import client
```

```
In [2]: mec = client.MultiEngineClient()
```

```
In [3]: mec.activate()
```

- ▶ Execute commands using %px

```
In [4]: %px from mpi4py import MPI
```

```
In [5]: %px print(MPI.Get_processor_name())
```

## Features – Interoperability

Good support for wrapping other MPI-based codes.

- ▶ You can use **Cython** (`cimport` statement).
- ▶ You can use **SWIG** (*typemaps* provided).
- ▶ You can use **F2Py** (`py2f()`/`f2py()` methods).
- ▶ You can use **Boost::Python** or **hand-written C** extensions.

**mpi4py** will allow you to use virtually any MPI based C/C++/Fortran code from Python.

# Hello World!

```
1  from mpi4py import MPI
2
3  rank = MPI.COMM_WORLD.Get_rank()
4  size = MPI.COMM_WORLD.Get_size()
5  name = MPI.Get_processor_name()
6
7  print ("Hello, World! "
8         "I am process %d of %d on %s" %
9         (rank, size, name))
```

# Hello World! – Wrapping with SWIG

## C source

```
1  /* file: helloworld.c */
2  void sayhello(MPI_Comm comm)
3  {
4      int size, rank;
5      MPI_Comm_size(comm, &size);
6      MPI_Comm_rank(comm, &rank);
7      printf("Hello, World! "
8             "I am process "
9             "%d of %d.\n",
10            rank, size);
11 }
```

## SWIG interface file

```
1  // file: helloworld.i
2  %module helloworld
3  %{
4      #include <mpi.h>
5      #include "helloworld.c"
6  }%
7
8  %include mpi4py/mpi4py.i
9  %mpi4py_typemap(Comm, MPI_Comm);
10
11 void sayhello(MPI_Comm comm);
```

## At the Python prompt ...

```
>>> from mpi4py import MPI
>>> import helloworld
>>> helloworld.sayhello(MPI.COMM_WORLD)
Hello, World! I am process 0 of 1.
```

# Hello World! – Wrapping with Boost.Python

```
1 // file: helloworld.cxx
2 #include <boost/python.hpp>
3 #include <mpi4py/mpi4py.h>
4 using namespace boost::python;
5
6 #include "helloworld.c"
7 static void wrap_sayhello(object py_comm) {
8     PyObject* py_obj = py_comm.ptr();
9     MPI_Comm *comm_p = PyMPIComm_Get(py_obj);
10    if (comm_p == NULL) throw_error_already_set();
11    sayhello(*comm_p);
12 }
13
14 BOOST_PYTHON_MODULE(helloworld) {
15     if (import_mpi4py() < 0) return;
16     def("sayhello", wrap_sayhello);
17 }
```

# Hello World! – Wrapping with F2Py

## Fortran 90 source

```
1  ! file: helloworld.f90
2  subroutine sayhello(comm)
3      use mpi
4      implicit none
5      integer :: comm, rank, size, ierr
6      call MPI_Comm_size(comm, size, ierr)
7      call MPI_Comm_rank(comm, rank, ierr)
8      print *, 'Hello, World! I am process ',rank,' of ',size,'.'
9  end subroutine sayhello
```

## At the Python prompt ...

```
>>> from mpi4py import MPI
>>> import helloworld
>>> fcomm = MPI.COMM_WORLD.py2f()
>>> helloworld.sayhello(fcomm)
Hello, World! I am process 0 of 1.
```



Overview

**Communicators**

Point to Point Communication

Collective Operations

Compute Pi

Mandelbrot Set

Dynamic Process Management

# Communicators

communicator = process group + communication context

## ▶ Predefined instances

- ▶ `COMM_WORLD`
- ▶ `COMM_SELF`
- ▶ `COMM_NULL`

## ▶ Accessors

- ▶ `rank = comm.Get_rank()` # or `comm.rank`
- ▶ `size = comm.Get_size()` # or `comm.size`
- ▶ `group = comm.Get_group()`

## ▶ Constructors

- ▶ `newcomm = comm.Dup()`
- ▶ `newcomm = comm.Create(group)`
- ▶ `newcomm = comm.Split(color, key)`

## Communicators – Create()

```
1  from mpi4py import MPI
2
3  comm = MPI.COMM_WORLD
4  group = comm.Get_group()
5
6  newgroup = group.Excl([0])
7  newcomm = comm.Create(newgroup)
8
9  if comm.rank == 0:
10     assert newcomm == MPI.COMM_NULL
11 else:
12     assert newcomm.size == comm.size - 1
13     assert newcomm.rank == comm.rank - 1
14
15 group.Free(); newgroup.Free()
16 if newcomm: newcomm.Free()
```

## Communicators – Split()

```
1  from mpi4py import MPI
2
3  world_rank = MPI.COMM_WORLD.Get_rank()
4  world_size = MPI.COMM_WORLD.Get_size()
5
6  if world_rank < world_size//2:
7      color = 55
8      key = -world_rank
9  else:
10     color = 77
11     key = +world_rank
12
13  newcomm = MPI.COMM_WORLD.Split(color, key)
14  # ...
15  newcomm.Free()
```

## Exercise #1

- a) Create a new process group containing the processes in the group of `COMM_WORLD` with **even** rank. Use the new group to create a new communicator.

**Tip:** use `Group.Incl()` or `Group.Range_incl()`

- b) Use `Comm.Split()` to split `COMM_WORLD` in two halves.
- ▶ The first half contains the processes with **even** rank in `COMM_WORLD`. The process rank ordering in the new communication is **ascending**.
  - ▶ The second half contains the processes with **odd** rank in `COMM_WORLD`. The process rank ordering in the new communication is **descending**.

Overview

Communicators

Point to Point Communication

Collective Operations

Compute Pi

Mandelbrot Set

Dynamic Process Management

▶ Blocking communication

▶ Python objects

```
comm.send(obj, dest=0, tag=0)  
obj = comm.recv(None, src=0, tag=0)
```

▶ Array data

```
comm.Send([array, count, datatype], dest=0, tag=0)  
comm.Recv([array, count, datatype], src=0, tag=0)
```

▶ Nonblocking communication

▶ Python objects

```
request = comm.isend(object, dest=0, tag=0}  
request.Wait()
```

▶ Array data

```
req1 = comm.Isend([array, count, datatype], dest=0, tag=0)  
req2 = comm.Irecv([array, count, datatype], src=0, tag=0)  
MPI.Request.Waitall([req1, req2])}
```

# PingPong

```
1 from mpi4py import MPI
2 comm = MPI.COMM_WORLD
3 assert comm.size == 2
4
5 if comm.rank == 0:
6     sendmsg = 777
7     comm.send(sendmsg, dest=1, tag=55)
8     recvmsg = comm.recv(source=1, tag=77)
9 else:
10    recvmsg = comm.recv(source=0, tag=55)
11    sendmsg = "abc"
12    comm.send(sendmsg, dest=0, tag=77)
```



# PingPing

```
1  from mpi4py import MPI
2  comm = MPI.COMM_WORLD
3  assert comm.size == 2
4
5  if comm.rank == 0:
6      sendmsg = 777
7      target = 1
8  else:
9      sendmsg = "abc"
10     target = 0
11
12     request = comm.isend(sendmsg, dest=target, tag=77)
13     recvmsg = comm.recv(source=target, tag=77)
14     request.Wait()
```

# Exchange

```
1  from mpi4py import MPI
2  comm = MPI.COMM_WORLD
3
4  sendmsg = [comm.rank]*3
5  right = (comm.rank + 1) % comm.size
6  left  = (comm.rank - 1) % comm.size
7
8  req1 = comm.isend(sendmsg, dest=right)
9  req2 = comm.isend(sendmsg, dest=left)
10 lmsg = comm.recv(source=left)
11 rmsg = comm.recv(source=right)
12
13 MPI.Request.Waitall([req1, req2])
14 assert lmsg == [left] * 3
15 assert rmsg == [right] * 3
```

## PingPing with NumPy arrays

```
1  from mpi4py import MPI
2  import numpy
3  comm = MPI.COMM_WORLD
4  assert comm.size == 2
5
6  if comm.rank == 0:
7      array1 = numpy.arange(10000, dtype='f')
8      array2 = numpy.empty(10000, dtype='f')
9      target = 1
10 else:
11     array1 = numpy.ones(10000, dtype='f')
12     array2 = numpy.empty(10000, dtype='f')
13     target = 0
14
15 request = comm.Isend([array1, MPI.FLOAT], dest=target)
16 comm.Recv([array2, MPI.FLOAT], source=target)
17 request.Wait()
```

## Exercise #2

- a) Modify *PingPong* example to communicate NumPy arrays.  
**Tip:** use `Comm.Send()` and `Comm.Recv()`
- b) Modify *Exchange* example to communicate NumPy arrays.  
Use nonblocking communication for both sending and receiving.  
**Tip:** use `Comm.Isend()` and `Comm.Irecv()`

Overview

Communicators

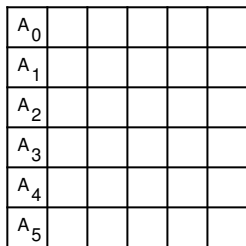
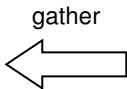
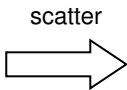
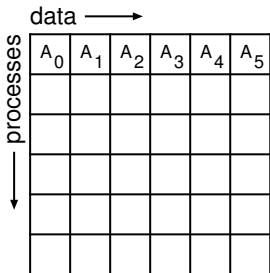
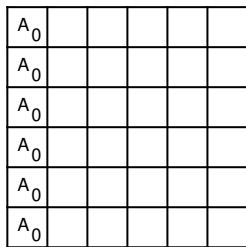
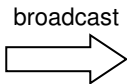
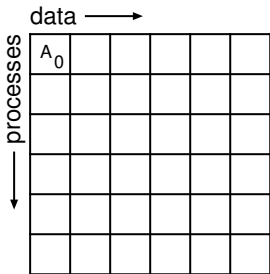
Point to Point Communication

**Collective Operations**

Compute Pi

Mandelbrot Set

Dynamic Process Management

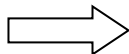


data →

processes ↓

A <sub>0</sub>					
B <sub>0</sub>					
C <sub>0</sub>					
D <sub>0</sub>					
E <sub>0</sub>					
F <sub>0</sub>					

allgather



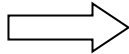
A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>

data →

processes ↓

A <sub>0</sub>	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
B <sub>0</sub>	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	B <sub>4</sub>	B <sub>5</sub>
C <sub>0</sub>	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	C <sub>4</sub>	C <sub>5</sub>
D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>	D <sub>4</sub>	D <sub>5</sub>
E <sub>0</sub>	E <sub>1</sub>	E <sub>2</sub>	E <sub>3</sub>	E <sub>4</sub>	E <sub>5</sub>
F <sub>0</sub>	F <sub>1</sub>	F <sub>2</sub>	F <sub>3</sub>	F <sub>4</sub>	F <sub>5</sub>

alltoall



A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>1</sub>	B <sub>1</sub>	C <sub>1</sub>	D <sub>1</sub>	E <sub>1</sub>	F <sub>1</sub>
A <sub>2</sub>	B <sub>2</sub>	C <sub>2</sub>	D <sub>2</sub>	E <sub>2</sub>	F <sub>2</sub>
A <sub>3</sub>	B <sub>3</sub>	C <sub>3</sub>	D <sub>3</sub>	E <sub>3</sub>	F <sub>3</sub>
A <sub>4</sub>	B <sub>4</sub>	C <sub>4</sub>	D <sub>4</sub>	E <sub>4</sub>	F <sub>4</sub>
A <sub>5</sub>	B <sub>5</sub>	C <sub>5</sub>	D <sub>5</sub>	E <sub>5</sub>	F <sub>5</sub>

# Broadcast

```
1 from mpi4py import MPI
2 comm = MPI.COMM_WORLD
3
4 if comm.rank == 0:
5     sendmsg = (7, "abc", [1.0,2+3j], {3:4})
6 else:
7     sendmsg = None
8
9 recvmsg = comm.bcast(sendmsg, root=0)
```



# Scatter

```
1 from mpi4py import MPI
2 comm = MPI.COMM_WORLD
3
4 if comm.rank == 0:
5     sendmsg = [i**2 for i in range(comm.size)]
6 else:
7     sendmsg = None
8
9 recvmsg = comm.scatter(sendmsg, root=0)
```

## Gather & Gather to All

```
1 from mpi4py import MPI
2 comm = MPI.COMM_WORLD
3
4 sendmsg = comm.rank**2
5
6 recvmsg1 = comm.gather(sendmsg, root=0)
7
8 recvmsg2 = comm.allgather(sendmsg)
```

## Reduce & Reduce to All

```
1 from mpi4py import MPI
2 comm = MPI.COMM_WORLD
3
4 sendmsg = comm.rank
5
6 recvmsg1 = comm.reduce(sendmsg, op=MPI.SUM, root=0)
7
8 recvmsg2 = comm.allreduce(sendmsg)
```

## Exercise #3

- a) Modify *Broadcast*, *Scatter*, and *Gather* example to communicate NumPy arrays.
- b) Write a routine implementing parallel *matrix–vector* product  $y = \text{matvec}(\text{comm}, A, x)$ .
  - ▶ the global matrix is dense and square.
  - ▶ matrix rows and vector entries have matching block distribution.
  - ▶ all processes own the same number of matrix rows.

**Tip:** use `Comm.Allgather()` and `numpy.dot()`

Overview

Communicators

Point to Point Communication

Collective Operations

**Compute Pi**

Mandelbrot Set

Dynamic Process Management

## Compute Pi

$$\pi = \int_0^1 \frac{4}{1+x^2} dx \approx \frac{1}{n} \sum_{i=0}^{n-1} \frac{4}{1 + \left(\frac{i+0.5}{n}\right)^2}$$

## Compute Pi – sequential

```
1  import math
2
3  def compute_pi(n):
4      h = 1.0 / n
5      s = 0.0
6      for i in range(n):
7          x = h * (i + 0.5)
8          s += 4.0 / (1.0 + x**2)
9      return s * h
10
11 n = 10
12 pi = compute_pi(n)
13 error = abs(pi - math.pi)
14
15 print ("pi is approximately %.16f, "
16       "error is %.16f" % (pi, error))
```

## Compute Pi – parallel [1]

```
1  from mpi4py import MPI
2  import math
3
4  def compute_pi(n, start=0, step=1):
5      h = 1.0 / n
6      s = 0.0
7      for i in range(start, n, step):
8          x = h * (i + 0.5)
9          s += 4.0 / (1.0 + x**2)
10     return s * h
11
12 comm = MPI.COMM_WORLD
13 nprocs = comm.Get_size()
14 myrank = comm.Get_rank()
```



## Compute Pi – parallel [2]

```
1  if myrank == 0:
2      n = 10
3  else:
4      n = None
5
6  n = comm.bcast(n, root=0)
7
8  mypi = compute_pi(n, myrank, nprocs)
9
10 pi = comm.reduce(mypi, op=MPI.SUM, root=0)
11
12 if myrank == 0:
13     error = abs(pi - math.pi)
14     print ("pi is approximately %.16f, "
15           "error is %.16f" % (pi, error))
```

## Exercise #4

Modify *Compute Pi* example to employ NumPy.

**Tip:** you can convert a Python `int/float` object to a NumPy *scalar* with `x = numpy.array(x)`.

Overview

Communicators

Point to Point Communication

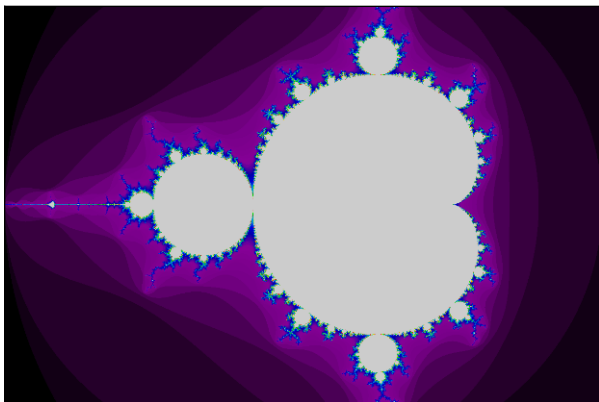
Collective Operations

Compute Pi

**Mandelbrot Set**

Dynamic Process Management

# Mandelbrot Set



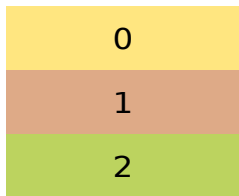
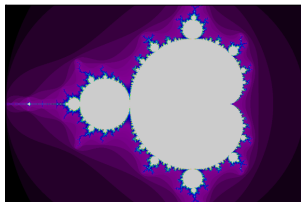
## Mandelbrot Set – sequential [1]

```
1  def mandelbrot(x, y, maxit):
2      c = x + y*1j
3      z = 0 + 0j
4      it = 0
5      while abs(z) < 2 and it < maxit:
6          z = z**2 + c
7          it += 1
8      return it
9
10 x1, x2 = -2.0, 1.0
11 y1, y2 = -1.0, 1.0
12 w, h = 150, 100
13 maxit = 127
```

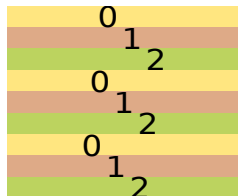
## Mandelbrot Set – sequential [2]

```
1 import numpy
2 C = numpy.zeros([h, w], dtype='i')
3 dx = (x2 - x1) / w
4 dy = (y2 - y1) / h
5 for i in range(h):
6     y = y1 + i * dy
7     for j in range(w):
8         x = x1 + j * dx
9         C[i, j] = mandelbrot(x, y, maxit)
10
11 from matplotlib import pyplot
12 pyplot.imshow(C, aspect='equal')
13 pyplot.spectral()
14 pyplot.show()
```

# Mandelbrot Set – partitioning



Block distribution



Cyclic distribution

## Mandelbrot Set – parallel, block [1]

```
1  def mandelbrot(x, y, maxit):
2      c = x + y*1j
3      z = 0 + 0j
4      it = 0
5      while abs(z) < 2 and it < maxit:
6          z = z**2 + c
7          it += 1
8      return it
9
10 x1, x2 = -2.0, 1.0
11 y1, y2 = -1.0, 1.0
12 w, h = 150, 100
13 maxit = 127
```



## Mandelbrot Set – parallel, block [2]

```
1  from mpi4py import MPI
2  import numpy
3
4  comm = MPI.COMM_WORLD
5  size = comm.Get_size()
6  rank = comm.Get_rank()
7
8  # number of rows to compute here
9  N = h // size + (h % size > rank)
10
11 # first row to compute here
12 start = comm.scan(N)-N
13
14 # array to store local result
15 Cl = numpy.zeros([N, w], dtype='i')
```

## Mandelbrot Set – parallel, block [3]

```
1  # compute owned rows
2
3  dx = (x2 - x1) / w
4  dy = (y2 - y1) / h
5  for i in range(N):
6      y = y1 + (i + start) * dy
7      for j in range(w):
8          x = x1 + j * dx
9          C1[i, j] = mandelbrot(x, y, maxit)
```

## Mandelbrot Set – parallel, block [4]

```
1  # gather results at root (process 0)
2
3  counts = comm.gather(N, root=0)
4  C = None
5  if rank == 0:
6      C = numpy.zeros([h, w], dtype='i')
7
8  rowtype = MPI.INT.Create_contiguous(w)
9  rowtype.Commit()
10
11 comm.Gatherv(sendbuf=[C1, MPI.INT],
12              recvbuf=[C, (counts, None), rowtype],
13              root=0)
14
15 rowtype.Free()
```

## Mandelbrot Set – parallel, block [5]

```
1  if comm.rank == 0:  
2      from matplotlib import pyplot  
3      pyplot.imshow(C, aspect='equal')  
4      pyplot.spectral()  
5      pyplot.show()
```

## Exercise #5

Measure the wall clock time  $T_i$  of local computations at each process for the *Mandelbrot Set* example with **block** and **cyclic** row distributions.

What row distribution is better regarding load balancing?

**Tip:** use `Wtime()` to measure wall time, compute the ratio  $T_{\max}/T_{\min}$  to compare load balancing.

Overview

Communicators

Point to Point Communication

Collective Operations

Compute Pi

Mandelbrot Set

**Dynamic Process Management**

# Dynamic Process Management

- ▶ Useful in assembling complex distributed applications. Can couple **independent parallel codes** written in **different languages**.
- ▶ Create new processes from a running program.
  - `Comm.Spawn()` and `Comm.Get_parent()`
- ▶ Connect two running applications together.
  - `Comm.Connect()` and `Comm.Accept()`

## Dynamic Process Management – Spawning

Spawning new processes is a *collective operation* that creates an **intercommunicator**.

- ▶ Local group is group of spawning processes (parent).
- ▶ Remote group is group of new processes (child).
- ▶ `Comm.Spawn()` lets parent processes spawn the child processes. It returns a new intercommunicator.
- ▶ `Comm.Get_parent()` lets child processes find intercommunicator to the parent group. Child processes have own `COMM_WORLD`.
- ▶ `Comm.Disconnect()` ends the parent–child connection. After that, both groups can continue running.



## Dynamic Process Management – Compute Pi (parent)

```
1  from mpi4py import MPI
2  import sys, numpy
3
4  comm = MPI.COMM_SELF.Spawn(sys.executable,
5                             args=['compute_pi-child.py'],
6                             maxprocs=5)
7
8  N = numpy.array(10, 'i')
9  comm.Bcast([N, MPI.INT], root=MPI.ROOT)
10 PI = numpy.array(0.0, 'd')
11 comm.Reduce(None, [PI, MPI.DOUBLE],
12             op=MPI.SUM, root=MPI.ROOT)
13 comm.Disconnect()
14
15 error = abs(PI - numpy.pi)
16 print ("pi is approximately %.16f, "
17        "error is %.16f" % (PI, error))
```

## Dynamic Process Management – Compute Pi (child)

```
1  from mpi4py import MPI
2  import numpy
3
4  comm = MPI.Comm.Get_parent()
5  size = comm.Get_size()
6  rank = comm.Get_rank()
7
8  N = numpy.array(0, dtype='i')
9  comm.Bcast([N, MPI.INT], root=0)
10 h = 1.0 / N; s = 0.0
11 for i in range(rank, N, size):
12     x = h * (i + 0.5)
13     s += 4.0 / (1.0 + x**2)
14 PI = numpy.array(s * h, dtype='d')
15 comm.Reduce([PI, MPI.DOUBLE], None,
16             op=MPI.SUM, root=0)
17
18 comm.Disconnect()
```

## Exercise #5

- a) Implement the *Compute Pi* **child** code in **C** or **C++** . Adjust the parent code in Python to spawn the new implementation.
- b) Compute and plot the *Mandelbrot Set* using spawning with parent/child codes implemented in Python.  
**Tip:** Reuse the provided parent code in Python and translate the child code in Fortran 90 to Python.

Do not hesitate to ask for help ...

- ▶ Mailing List: [mpi4py@googlegroups.com](mailto:mpi4py@googlegroups.com)
- ▶ Mail&Chat: [dalcin1@gmail.com](mailto:dalcin1@gmail.com)

# Thanks!