# GPU Programming using Python
## (PyCUDA and PyOpenCL)

Ki-Hwan Kim

Department of Physics, Korea University

2010/5/28

## Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

## Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

# Scientific computing environments

## Matlab, Mathematica, Maple, ...

- simple and clean syntax
- high productivity
- tight integration of calculation and visualization

=> Do not work for some problems, at least not in an easy way

## Traditional languages (Fortran, C, C++, JAVA)

- flexible and versatile
- high performance

=> Complicated work, Low productivity

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

# Scientific computing environments

### Matlab, Mathematica, Maple, ...

- simple and clean syntax
- high productivity
- tight integration of calculation and visualization

=> Do not work for some problems, at least not in an easy way

### Traditional languages (Fortran, C, C++, JAVA)

- flexible and versatile
- high performance

=> Complicated work, Low productivity

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

# Scientific computing environments

## Matlab, Mathematica, Maple, ...

- simple and clean syntax
- high productivity
- tight integration of calculation and visualization

=> Do not work for some problems, at least not in an easy way

## Traditional languages (Fortran, C, C++, JAVA)

- flexible and versatile
- high performance

=> Complicated work, Low productivity

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

# Scientific computing environments

## Matlab, Mathematica, Maple, ...

- simple and clean syntax
- high productivity
- tight integration of calculation and visualization

=> Do not work for some problems, at least not in an easy way

## Traditional languages (Fortran, C, C++, JAVA)

- flexible and versatile
- high performance

=> Complicated work, Low productivity

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

# Scientific computing environments (Cont.)

## Scripting using Python

- simple and clean syntax
- gluing other favorate simulation, visualization, data analysis
- high productivity
- high performance (gluing Fortran, C)
- Easy integration with GPU computing (PyCUDA, PyOpenCL)

=> Build my own Matlab-like scientific computing environment, tailored to my specific needs with high performance

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Scientific computing environments

# Scientific computing environments (Cont.)

## Scripting using Python

- simple and clean syntax
- gluing other favorate simulation, visualization, data analysis
- high productivity
- high performance (gluing Fortran, C)
- Easy integration with GPU computing (PyCUDA, PyOpenCL)

=> Build my own Matlab-like scientific computing environment,
tailored to my specific needs with high performance

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

**Formula**
Python, PyCUDA, PyOpenCL
PyCUDA vs PyOpenCL

# Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA, PyOpenCL
PyCUDA vs PyOpenCL

# Mathematical Formula

Gaussian and Sin functions

$$\begin{cases} a(x) & = e^{-\frac{x^2}{2}} \\ b(x) & = \sin(5x) \end{cases}$$

$$\Rightarrow c(x) = a(x) \times b(x) = e^{-\frac{x^2}{2}} \sin(5x)$$

Discretize Condition

$$-5 \leq x < 5 \quad and \quad \Delta_x = 0.01$$

$$\rightarrow n = 1000$$

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

**Formula**
Python, PyCUDA, PyOpenCL
PyCUDA vs PyOpenCL

## Mathematical Formula

Gaussian and Sin functions

$$\begin{cases} a(x) & = e^{-\frac{x^2}{2}} \\ b(x) & = \sin{(5x)} \end{cases}$$

$$\Rightarrow c(x) = a(x) \times b(x) = e^{-\frac{x^2}{2}} \sin(5x)$$

Discretize Condition

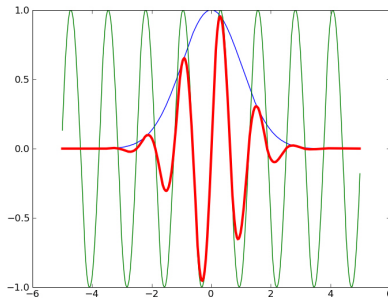$$-5 \leq x < 5 \quad and \quad \Delta_x = 0.01$$

$$\rightarrow n = 1000$$

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA, PyOpenCL
PyCUDA vs PyOpenCL

## Outline

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA, PyOpenCL**
PyCUDA vs PyOpenCL

# Python

```
1  #!/usr/bin/env python
2  import numpy as np
3
4  x = np.arange(-5,5,0.01,'f')
5  a = np.exp(-(x**2)/2)
6  b = np.sin(5*x)
7
8  c = a*b
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA, PyOpenCL**
PyCUDA vs PyOpenCL

# Python (Plot the graph)

```
1  import matplotlib.pyplot as pl
2
3  pl.plot(x,a,x,b)
4  pl.plot(x,c,linewidth=3)
5  pl.savefig('./pics/mul_array.png')
6  pl.show()
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA, PyOpenCL**
PyCUDA vs PyOpenCL

# PyCUDA

```python
1   #!/usr/bin/env python
2   import numpy as np
3
4   x = np.arange(-5,5,0.01,'f')
5   a = np.exp(-(x**2)/2)
6   b = np.sin(5*x)
7
8   import pycuda.driver as cuda
9   import pycuda.autoinit
10
11  a_gpu = cuda.to_device(a)
12  b_gpu = cuda.to_device(b)
13  c_gpu = cuda.mem_alloc(a.nbytes)
14
15  gpu_mul(np.int32(a.size), a_gpu, b_gpu, c_gpu,
16          block=(256,1,1),grid=(4,1))
17
18  c = np.zeros_like(a)
19  cuda.memcpy_dtoh(c, c_gpu)
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA, PyOpenCL
PyCUDA vs PyOpenCL

# PyCUDA (Cont.)

```
1   kernels = """
2     __global__ void multiply(int nx, float *a,
3           float *b, float *c){
4
5       int idx = threadIdx.x;
6
7       if(idx<nx) c[idx] = a[idx]*b[idx];
8     }
9   """
10
11  from pycuda.compiler import SourceModule
12  mod = SourceModule(kernels)
13  gpu_mul = mod.get_function("multiply")
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA, PyOpenCL**
PyCUDA vs PyOpenCL

# PyOpenCL

```python
1   #!/usr/bin/env python
2   import numpy as np
3
4   x = np.arange(-5,5,0.01,'f')
5   a = np.exp(-(x**2)/2)
6   b = np.sin(5*x)
7
8   import pyopencl as cl
9   ctx = cl.create_some_context()
10  queue = cl.CommandQueue(ctx)
11
12  mf=cl.mem_flags
13  a_gpu = cl.Buffer(ctx, mf.COPY_HOST_PTR, hostbuf=a)
14  b_gpu = cl.Buffer(ctx, mf.COPY_HOST_PTR, hostbuf=b)
15  c_gpu = cl.Buffer(ctx, mf.COPY_HOST_PTR, size=a.nbytes)
16
17  gpu_mul(queue, (a.size,) np.int32(a.size), a_gpu, b_gpu, c_gpu)
18
19  c = np.zeros_like(a)
20  cl.enqueue_read_buffer(queue, c_gpu, c)
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA, PyOpenCL**
PyCUDA vs PyOpenCL

# PyOpenCL (Cont.)

```
1   kernels = """
2     __kernel void multiply(int nx, __global float *a,
3             __global float *b, __global float *c) {
4
5       int idx = get_global_id(0);
6
7       if (idx<nx) c[idx] = a[idx]*b[idx];
8     }
9   """
10
11  prg = cl.Program(ctx, kernels).build()
12  gpu_mul = prg.multiply
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA, PyOpenCL
**PyCUDA vs PyOpenCL**

# Outline

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA, PyOpenCL
**PyCUDA vs PyOpenCL**

# PyCUDA vs PyOpenCL

```
1   import pycuda.driver as cuda
2   import pycuda.autoinit
3
4   a_gpu = cuda.to_device(a)
5   b_gpu = cuda.to_device(b)
6   c_gpu = cuda.mem_alloc(a.nbytes)
7   gpu_mul(np.int32(a.size), a_gpu, b_gpu, c_gpu,
8           block=(256,1,1),grid=(4,1))
9   cuda.memcpy_dtoh(c, c_gpu)
```

```
1    import pyopencl as cl
2    ctx = cl.create_some_context()
3    queue = cl.CommandQueue(ctx)
4
5    mf=cl.mem_flags
6    a_gpu = cl.Buffer(ctx, mf.COPY_HOST_PTR, hostbuf=a)
7    b_gpu = cl.Buffer(ctx, mf.COPY_HOST_PTR, hostbuf=b)
8    c_gpu = cl.Buffer(ctx, mf.COPY_HOST_PTR, size=a.nbytes)
9    gpu_mul(queue, (a.size,) np.int32(a.size), a_gpu, b_gpu, c_gpu)
10   cl.enqueue_read_buffer(queue, c_gpu, c)
```

Why Python for scientific computations?
**Example 1) 1-D Array Multiplication**
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA, PyOpenCL
**PyCUDA vs PyOpenCL**

# PyCUDA vs PyOpenCL (Cont.)

```
1  kernels = """
2    __global__ void multiply(int nx, float *a,...) {
3      int idx = threadIdx.x;
4      if(idx<nx) c[idx] = a[idx]*b[idx];
5    }
6  """
7
8  from pycuda.compiler import SourceModule
9  mod = SourceModule(kernels)
10 gpu_mul = mod.get_function("multiply")
```

```
1  kernels = """
2    __kernel void multiply(int nx, __global float *a,...) {
3      int idx = get_global_id(0);
4      if(idx<nx) c[idx] = a[idx]*b[idx];
5    }
6  """
7
8  prg = cl.Program(ctx, kernels).build()
9  gpu_mul = prg.multiply
```

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

## Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

## Mathematical Formula

**2-D Wave Equation**

$$\frac{\partial^2 f}{\partial t^2} = c^2 \nabla^2 f$$

$$\rightarrow \frac{\partial^2 f}{\partial t^2} = c^2 \left( \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} \right)$$

**Simple Finite-Difference Scheme**

$$\frac{f_{i,j}^{n+1} - 2f_{i,j}^n + f_{i,j}^{n-1}}{\Delta_t^2} = c_{i,j}^2 \left( \frac{f_{i+1,j}^n - 2f_{i,j}^n + f_{i-1,j}^n}{\Delta_x^2} + \frac{f_{i,j+1}^n - 2f_{i,j}^n + f_{i,j-1}^n}{\Delta_y^2} \right)$$

$$f_{i,j}^{n+1} = \tilde{c}_{i,j}^2 \left( f_{i+1,j}^n + f_{i-1,j}^n + f_{i,j+1}^n + f_{i,j-1}^n - 4f_{i,j}^n \right) + 2f_{i,j}^n - f_{i,j}^{n-1}$$

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

## Mathematical Formula

2-D Wave Equation

$$\frac{\partial^2 f}{\partial t^2} = c^2 \nabla^2 f$$

$$\rightarrow \frac{\partial^2 f}{\partial t^2} = c^2 \left( \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} \right)$$

Simple Finite-Difference Scheme

$$\frac{f_{i,j}^{n+1} - 2f_{i,j}^n + f_{i,j}^{n-1}}{\Delta_t^2} = c_{i,j}^2 \left( \frac{f_{i+1,j}^n - 2f_{i,j}^n + f_{i-1,j}^n}{\Delta_x^2} + \frac{f_{i,j+1}^n - 2f_{i,j}^n + f_{i,j-1}^n}{\Delta_y^2} \right)$$

$$f_{i,j}^{n+1} = \tilde{c}_{i,j}^2 \left( f_{i+1,j}^n + f_{i-1,j}^n + f_{i,j+1}^n + f_{i,j-1}^n - 4f_{i,j}^n \right) + 2f_{i,j}^n - f_{i,j}^{n-1}$$

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
**Python, PyCUDA**
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA**
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Python

```
1   import numpy as np
2
3   c = np.ones((1000,1000),'f')*0.25
4   f = np.zeros_like(c)
5   g = np.zeros_like(c)
6
7   ii = (slice(1,-1), slice(1,-1))
8   for tn in xrange(1000):
9       g[400,500] += np.sin(0.1*tn)
10      f[ii] = c[ii]*(g[2:,1:-1]+g[:-2,1:-1]+g[1:-1,2:]+g[1:-1,:-2]
11                     -4g[ii])+2g[ii]-f[ii]
12      g[ii] = c[ii]*(f[2:,1:-1]+f[:-2,1:-1]+f[1:-1,2:]+f[1:-1,:-2]
13                     -4f[ii])+2f[ii]-g[ii]
```

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
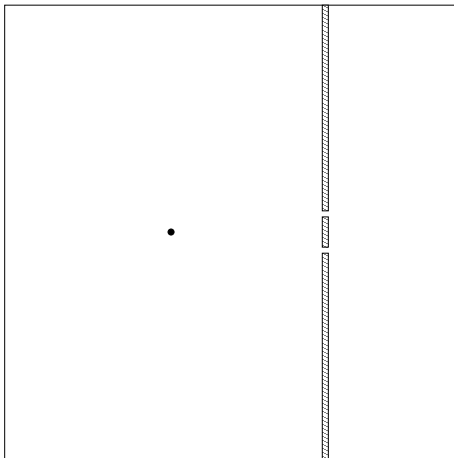Utilize Multi-GPUs with MPI

# PyCUDA

```
1   kernels = """
2     __global__ void update_src(int idx, int tn, float *g) {
3       g[idx] += sin(0.1*tn);
4     }
5
6     __global__ void update(int nx, int ny, float *c, float *f,
            float *g) {
7       int idx = blockIdx.x*blockDim.x + threadIdx.x;
8       int i = idx/ny, j = idx%ny;
9
10      if(i>0 && j>0 && i<nx−1 && j<ny−1)
11        f[idx] = c[idx]*(g[idx+ny]+g[idx−ny]+g[idx+1]+g[idx−1]
12              −4*g[idx])+2*g[idx]−f[idx];
13    }
14  """
15
16  from pycuda.compiler import SourceModule
17  mod = SourceModule(kernels)
18  update_src = mod.get_function("update_src")
19  update = mod.get_function("update")
```
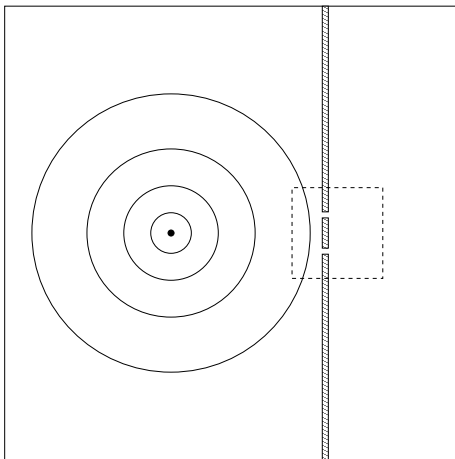
Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
**Python, PyCUDA**
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# PyCUDA (Cont.)

```
1   import pycuda.driver as cuda
2   import pycuda.autoinit
3
4   c = np.ones((nx,ny),'f')*0.25
5   f = np.zeros_like(c)
6   c_gpu = cuda.to_device(c)
7   f_gpu = cuda.to_device(f)
8   g_gpu = cuda.to_device(f)
9   ...
10  cuda.memcpy_htod(c_gpu,c)
11
12  Db, Dg = (256,1,1), (nx*ny/256+1,1)
13  nnx, nny = np.int32(nx), np.int32(ny)
14  src_pt = np.int32(nx/2*ny + ny/2)
15  ...
16  for tn in xrange(10000):
17      update_src(src_pt,np.int32(tn),f_gpu,(1,1,1),(1,1))
18      update(nnx,nny,c_gpu,f_gpu,g_gpu,Db,Dg)
19      update(nnx,nny,c_gpu,g_gpu,f_gpu,Db,Dg)
20      if tn%100==0: cuda.memcpy_dtoh(f,f_gpu)
```
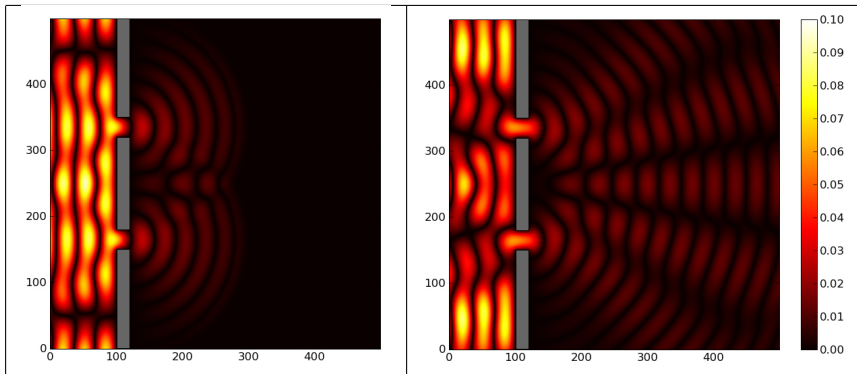
Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
**Python, PyCUDA**
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Example - Double Slit Interference

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Double Slit Interference (source point)

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Double Slit Interference (result)

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
Python, PyCUDA
**Shared Memory Optimization**
Utilize Multi-GPUs with MPI

## Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
**Shared Memory Optimization**
Utilize Multi-GPUs with MPI

# Shared Memory Optimize

```
1  __global__ void update(int nx, int ny, float *c, float *f,
       float *g) {
2      int idx = blockIdx.x*blockDim.x + threadIdx.x;
3      int i = idx/ny, j = idx%ny;
4
5      if(i>0 && j>0 && i<nx-1 && j<ny-1)
6          f[idx] = c[idx]*(g[idx+ny]+g[idx-ny]
7              +g[idx+1]+g[idx-1]-4*g[idx])+2*g[idx]-f[idx];
8  }
```

Duplicated arrays in a thread block
=> transfer to shared memory, and reuse
=> reduce global memory access
Target: g[idx+1], g[idx-1], g[idx], g[idx]

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
**Shared Memory Optimization**
Utilize Multi-GPUs with MPI

## Shared Memory Optimize (Cont.)

```
1   __global__ void update(int nx, int ny, float *c, float *f,
         float *g) {
2       int tx = threadIdx.x;
3       int idx = blockIdx.x*blockDim.x + tx;
4
5       extern __shared__ float gs[];
6       gs[tx+1] = g[idx];
7       int i = idx/ny, j = idx%ny;
8       if(j>0 && j<ny-1) {
9           if(tx==0) gs[tx] = g[idx-1];
10          if(tx==blockDim.x-1) gs[tx+2] = g[idx+1];
11      }
12      __syncthreads()
13
14      if(i>0 && j>0 && i<nx-1 && j<ny-1)
15          f[idx] = c[idx]*(g[idx+ny]+g[idx-ny]
16              +gs[tx+2]+gs[tx]-4*gs[tx+1])+2*gs[tx+1]-f[idx];
17  }
```

```
1   update(nnx,nny,c_gpu,f_gpu,g_gpu,Db,Dg,shared=(256+2)*4)
2   update(nnx,nny,c_gpu,g_gpu,f_gpu,Db,Dg,shared=(256+2)*4)
```

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

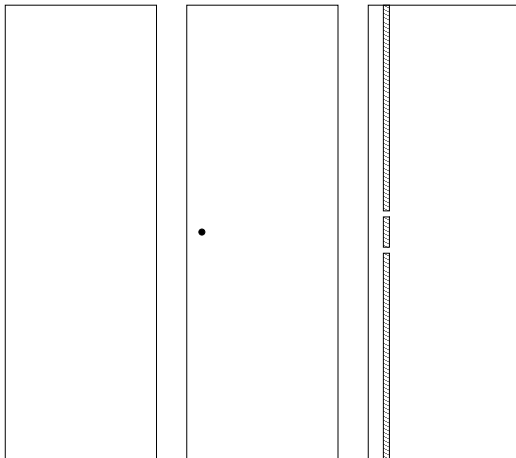# Shared Memory Optimize (Cont.)

```
1   __global__ void update(int nx, int ny, float *c, float *f,
        float *g) {
2       int tx = threadIdx.x;
3       int idx = blockIdx.x*blockDim.x + tx;
4
5       extern __shared__ float gs[];
6       gs[tx+1] = g[idx];
7       int i = idx/ny, j = idx%ny;
8       if(j>0 && j<ny-1) {
9           if(tx==0) gs[tx] = g[idx-1];
10          if(tx==blockDim.x-1) gs[tx+2] = g[idx+1];
11      }
12      __syncthreads()
13
14      if(i>0 && j>0 && i<nx-1 && j<ny-1)
15          f[idx] = c[idx]*(g[idx+ny]+g[idx-ny]
16              +gs[tx+2]+gs[tx]-4*gs[tx+1])+2*gs[tx+1]-f[idx];
17  }
```

```
1   update(nnx,nny,c_gpu,f_gpu,g_gpu,Db,Dg,shared=(256+2)*4)
2   update(nnx,nny,c_gpu,f_gpu,g_gpu,Db,Dg,shared=(256+2)*4)
```

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
Python, PyCUDA
Shared Memory Optimization
**Utilize Multi-GPUs with MPI**

# Outline

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Divide Domain

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Move to MPI

```python
import pycuda.driver as cuda
import pycuda.autoinit

c = sc.ones((nx,ny),'f')*0.25
f = sc.zeros_like(c)
c_gpu = cuda.to_device(c)
f_gpu = cuda.to_device(f)
g_gpu = cuda.to_device(f)

# set the c array with geometry
cuda.memcpy_htod(c_gpu,c)

Db, Dg = (256,1,1), (nx*ny/256+1,1)
nnx, nny = sc.int32(nx), sc.int32(ny)
src_pt = sc.int32(nx/2*ny + ny/2)

# SourceModule

for tn in xrange(10000):
    update_src(src_pt,...)
    update(...,f_gpu,g_gpu,...)
    update(...,g_gpu,f_gpu,...)

    if tn%100==0:
        cuda.memcpy_dtoh(f,f_gpu)
```

```python
import pycuda.driver as cuda
import boostmpi as mpi

cuda.init()
dev = cuda.Device(mpi.rank)
ctx = dev.make_context()

# memory allocate

if mpi.rank == 2:
    # set the c array with geometry
    cuda.memcpy_htod(c_gpu,c)

Db, Dg = ...
nnx, nny = ...
if mpi.rank == 1:
    src_pt = ...

# SourceModule

for tn in xrange(10000):
    if mpi.rank == 1:
        update_src(src_pt,...)
    update(...,f_gpu,g_gpu,...)
    update(...,g_gpu,f_gpu,...)

ctx.pop()
```
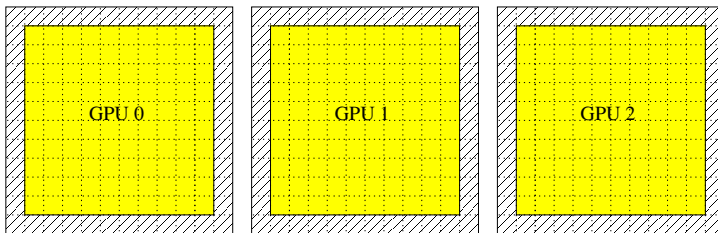
Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Gather data

```
1   nxx,  nyy  =  nx−2,  ny−2
2
3   if  mpi.rank  ==  0:
4       output  =  sc.zeros(3*nxx,  nyy,  'f')
5
6   for  tn  in  xrange(1000):
7         ...
8         if  tn  >  1000:
9             if  mpi.rank  ==  0:
10                cuda.memcpy_dtoh(f,  f_gpu)
11                output[:nxx,:]  =  f[1:−1,1:−1]
12                output[nxx:2*nxx,:]  =  mpi.world.recv(1,3)
13                output[2*nxx:3*nxx,:]  =  mpi.world.recv(2,3)
14            else:
15                cuda.memcpy_dtoh(f,  f_gpu)
16                mpi.world.send(0,  3,  f[1:−1,1:−1]
```
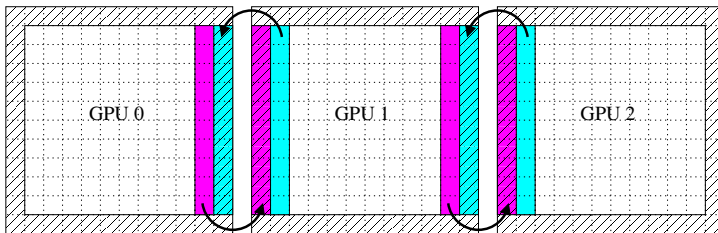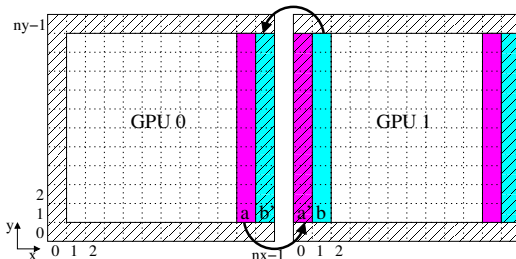
Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Exchange Boundaries

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
**Example 2) 2-D Wave Equation**

Formula
Python, PyCUDA
Shared Memory Optimization
**Utilize Multi-GPUs with MPI**

# Memcpy Offset



## offset

a: $\text{int(a\_gpu)} + ((nx-2)*ny+1)*sof$

a': $\text{int(a\_gpu)} + 1*sof$

b: $\text{int(a\_gpu)} + (ny+1)*sof$

b': $\text{int(a\_gpu)} + ((nx-1)*ny+1)*sof$

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# Exchange Boundaries

```
1   nbof = np.nbytes['float32']
2   dtof = np.dtype('float32')
3
4   def send(rank, tag_mark, nx, ny, a_gpu):
5       if mpi.rank < rank:
6           offset = int(a_gpu)+((nx-2)*ny+1)*nbof
7       elif mpi.rank > rank:
8           offset = int(a_gpu)+(ny+1)*nbof
9       mpi.world.send(rank, tag_mark, \
10              cuda.from_device(offset, (ny-2,), dtof) )
11
12  def recv(rank, tag_mark, nx, ny, a_gpu):
13      if mpi.rank > rank:
14          offset = int(a_gpu)+1*nbof
15      elif mpi.rank < rank:
16          offset = int(a_gpu)+((nx-1)*ny+1)*nbof
17      cuda.memcpy_htod(offset, mpi.world.recv(rank, tag_mark))
```

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

## Exchange Boundaries (Cont.)

```
1   def exchange(nx, ny, a_gpu):
2       if mpi.rank == 0:
3           send(1, 0, nx, ny, a_gpu)
4           recv(1, 0, nx, ny, a_gpu)
5       if mpi.rank == 1:
6           recv(0, 0, nx, ny, a_gpu)
7           send(0, 1, nx, ny, a_gpu)
8           send(2, 0, nx, ny, a_gpu)
9           recv(2, 1, nx, ny, a_gpu)
10      if mpi.rank == 2:
11          recv(1, 0, nx, ny, a_gpu)
12          send(1, 1, nx, ny, a_gpu)
```

```
1   for tn in xrange(10000):
2       if mpi.rank == 1:
3           update_src(src_pt, ...)
4       update(..., f_gpu, g_gpu, ...)
5       exchange(nx, ny, f_gpu)
6       update(..., g_gpu, f_gpu, ...)
7       exchange(nx, ny, g_gpu)
```

Why Python for scientific computations?
Example 1) 1-D Array Multiplication
Example 2) 2-D Wave Equation

Formula
Python, PyCUDA
Shared Memory Optimization
Utilize Multi-GPUs with MPI

# *Have your Fun with Python!*

## THANK YOU