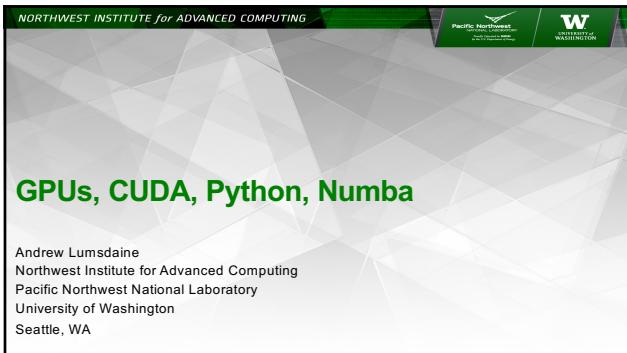
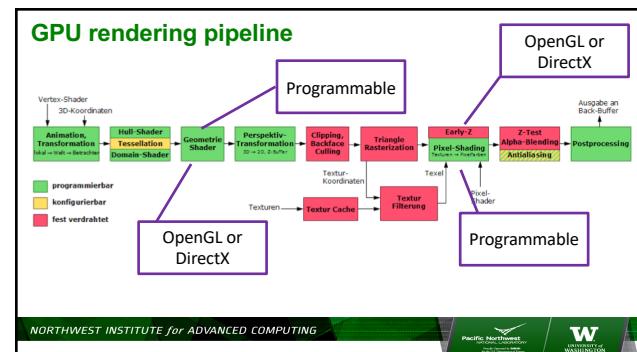
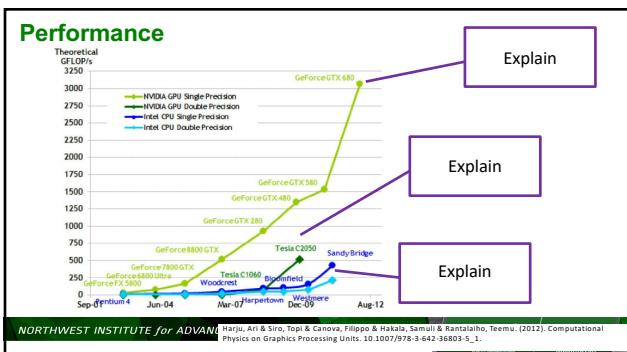
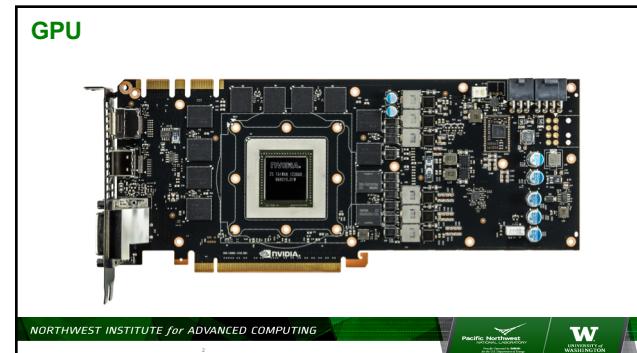


NORTHWEST INSTITUTE for ADVANCED COMPUTING



GPUs, CUDA, Python, Numba

Andrew Lumsdaine
Northwest Institute for Advanced Computing
Pacific Northwest National Laboratory
University of Washington
Seattle, WA



Compute = Drawing, Kernel = Shader

```

float saxpy (
    float2 coords : TEXCOORD0,
    uniform sampler2D textureY,
    uniform sampler2D textureX,
    uniform float alpha ) : COLOR
{
    float result;
    float y = tex2D(textureY,coords);
    float x = tex2D(textureX,coords);
    result = y + alpha * x;
    return result;
}

for (int i=0; i<N; i++)
    dataY[i] = dataY[i] + alpha * dataX[i];
    
```

Lots and lots of graphics code

```

#include <stdio.h>
#include <stdlib.h>
#include <GL/glew.h>
#include <GL/glut.h>

int main(int argc, char **argv) {
    // declare texture size, the actual data will be a vector
    // of size texSize*texSize*4
    int texSize = 2;
    // create texture
    float* data = (float*)malloc(sizeof(float)*texSize*texSize*4*texSize*texSize*4);
    float* result = (float*)malloc(sizeof(float)*texSize*texSize*4*texSize*texSize*4);
    for (int i=0; i<texSize*texSize*4; i++)
        data[i] = i+1.0;
    // set up glut to get valid GL context and
    // get extension entry points
    glutInit (&argc, argv);
    glutCreateWindow("TEST1");
    glewInit();
    
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING | Pacific Northwest National Laboratory | University of Washington

Lots and lots

```
// viewport transform for 1:1 pixel=texel=data mapping
glMatrixMode(GL_PROJECTION);
glLoadIdentity();
glOrtho(0, texSize, 0, texSize);
glMatrixMode(GL_MODELVIEW);
glLoadIdentity();
glViewport(0,0,texSize,texSize);
// create FBO and bind it (that is, use offscreen render target)
GLuint fb;
glGenFramebuffersEXT(1,&fb);
glBindFramebufferEXT(GL_FRAMEBUFFER_EXT,fb);
// create texture
GLuint tex;
glGenTextures (1,&tex);
glBindTexture(GL_TEXTURE_RECTANGLE_ARB, tex);
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest National Laboratory

UNIVERSITY OF WASHINGTON

Still more

```
// set texture parameters
glTexParameteri(GL_TEXTURE_RECTANGLE_ARB,
                GL_TEXTURE_MIN_FILTER, GL_NEAREST);
glTexParameteri(GL_TEXTURE_RECTANGLE_ARB,
                GL_TEXTURE_MAG_FILTER, GL_NEAREST);
glTexParameteri(GL_TEXTURE_RECTANGLE_ARB,
                GL_TEXTURE_WRAP_S, GL_CLAMP);
glTexParameteri(GL_TEXTURE_RECTANGLE_ARB,
                GL_TEXTURE_WRAP_T, GL_CLAMP);
// define texture with floating point format
glTexImage2D(GL_TEXTURE_RECTANGLE_ARB,0,GL_RGBA32F_ARB,
             texSize,texSize,0,GL_RGBA,GL_FLOAT,0);
// attach texture
glFramebufferTexture2DEXT(GL_FRAMEBUFFER_EXT,
                         GL_COLOR_ATTACHMENT0_EXT,
                         GL_TEXTURE_RECTANGLE_ARB, tex,0);
```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest National Laboratory

UNIVERSITY OF WASHINGTON

And.... Done

```
// transfer data to texture
glTexSubImage2D(GL_TEXTURE_RECTANGLE_ARB,0,0,0,texSize,texSize,
                GL_RGBA,GL_FLOAT,data);
// and read back
glReadBuffer(GL_COLOR_ATTACHMENT0_EXT);
glReadPixels(0, 0, texSize, texSize,GL_RGBA,GL_FLOAT,result);
printf("Data before roundtrip:\n");
for (int i=0; i<texSize*texSize*4; i++)
    printf("%c\n",data[i]);
printf("Data after roundtrip:\n");
for (int i=0; i<texSize*texSize*4; i++)
    printf("%c\n",result[i]);
// free data
free(data);
free(result);
glDeleteFramebuffersEXT (1,&fb);
glDeleteTextures (1,&tex);
return 0;
```

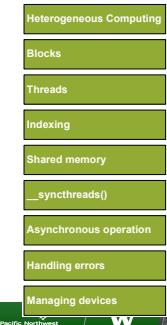
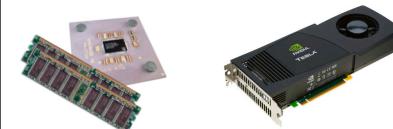
NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest National Laboratory

UNIVERSITY OF WASHINGTON

GPU programming

- Host: CPU and its memory
 - Host memory
- Device: GPU and its memory
 - Device memory

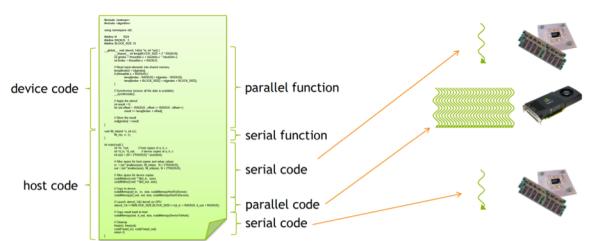


NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest National Laboratory

UNIVERSITY OF WASHINGTON

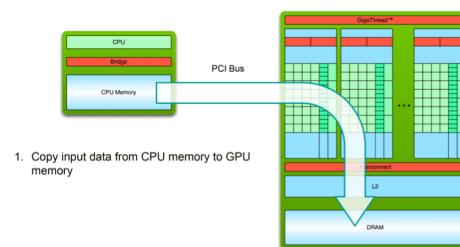
Heterogeneous computing



NORTHWEST INSTITUTE for ADVANCED COMPUTING

© NVIDIA Corporation 2011

Simple processing flow

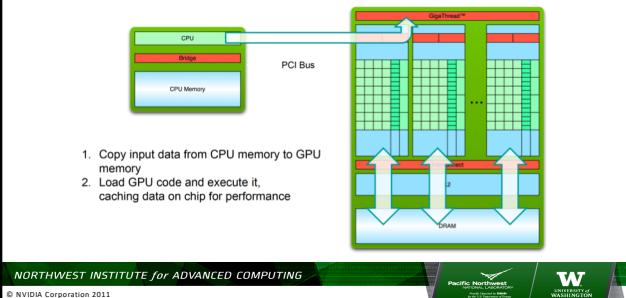


NORTHWEST INSTITUTE for ADVANCED COMPUTING

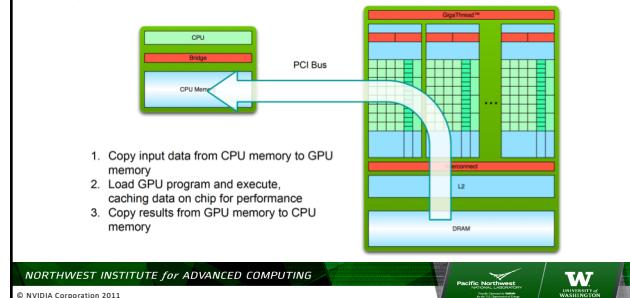
Pacific Northwest National Laboratory

UNIVERSITY OF WASHINGTON

Simple processing flow



Simple processing flow

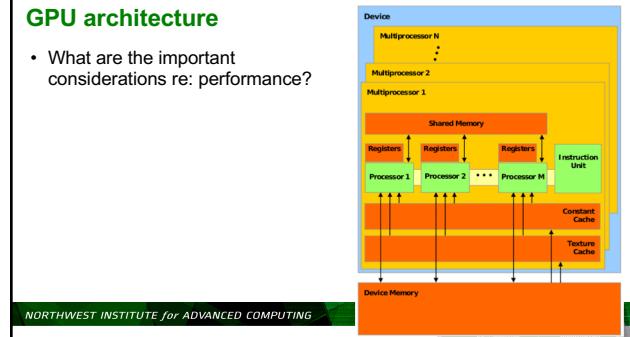


GPU



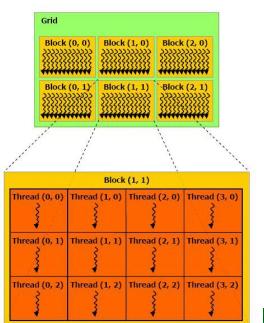
GPU architecture

- What are the important considerations re: performance?



Grids and blocks

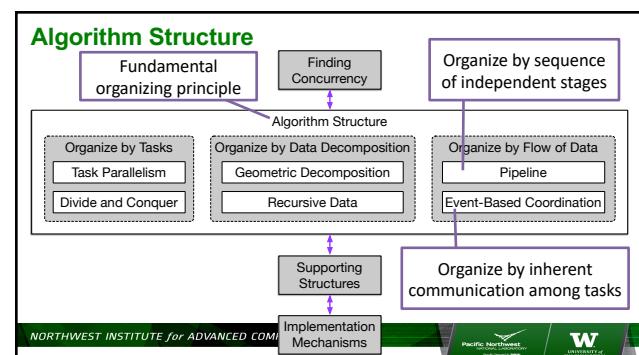
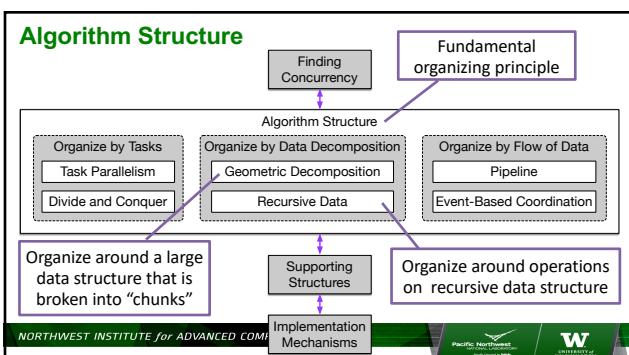
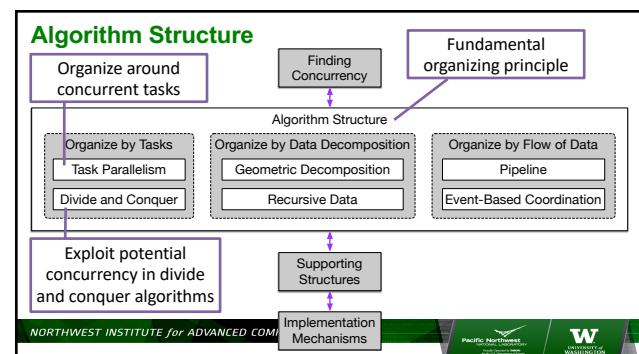
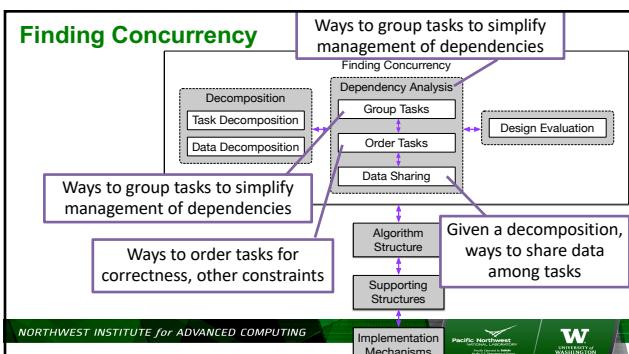
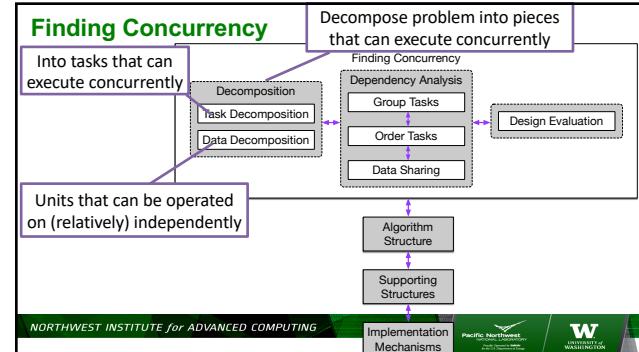
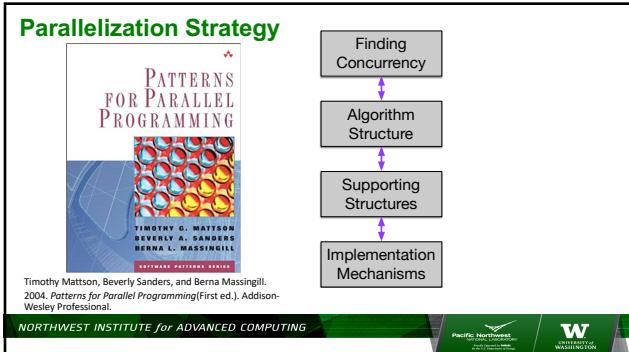
- Thread - Distributed by the CUDA runtime (threadIdx)
- Block - A user defined group of 1 to ~512 threads (blockIdx)
- Grid - A group of one or more blocks. A grid is created for each CUDA kernel function called

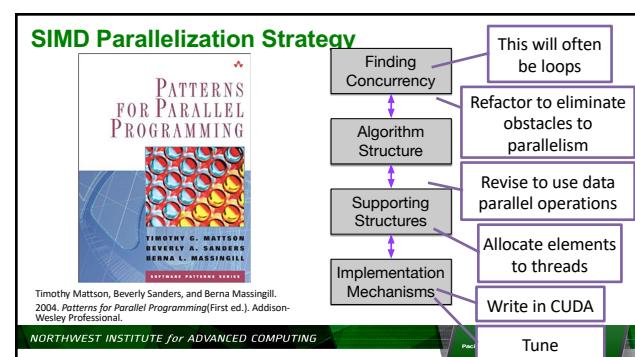
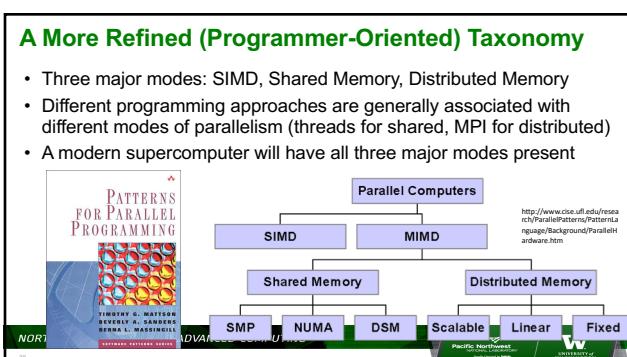
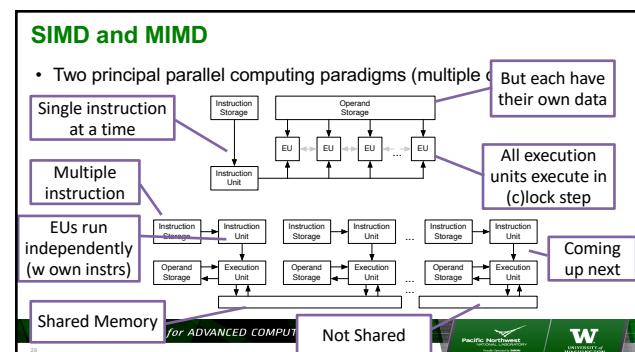
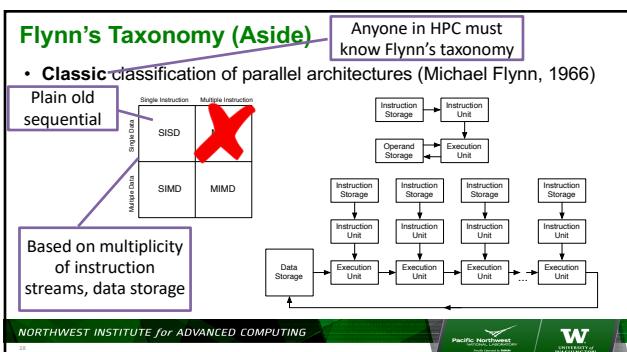
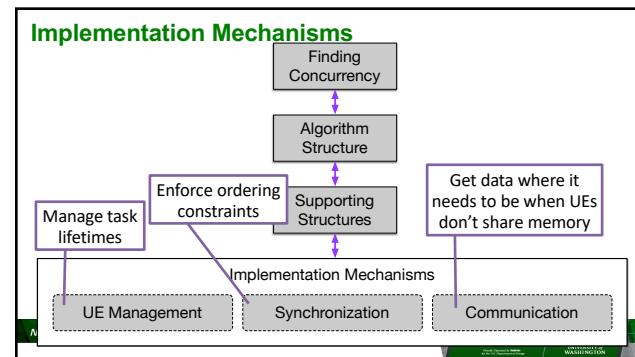
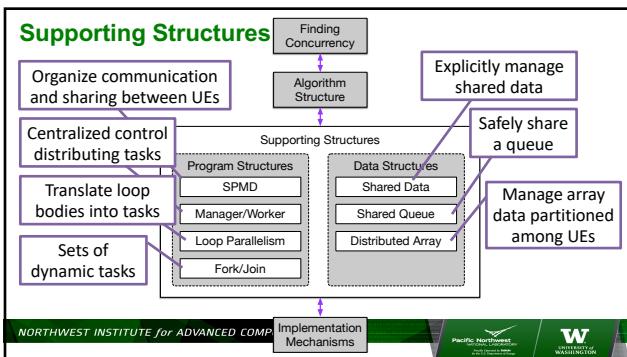


GPU Programming

- Transfer data from CPU memory to GPU memory
- Compute in GPU
- Transfer data back
- Computation is done with huge number of identical processing units
- Each PU executes exactly same kernel
- Partition the data among the processors – use thread idx, block idx, and block dimension to compute a unique (global)
- Use global index in the kernel to determine data each thread should work on





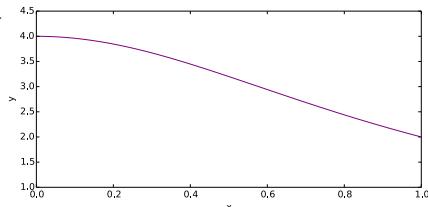


Example

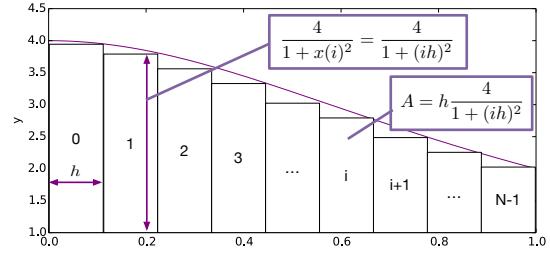
- Find the value of π

- Using formula

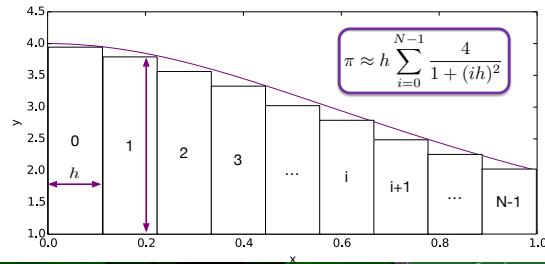
$$\pi = \int_0^1 \frac{4}{1+x^2} dx$$



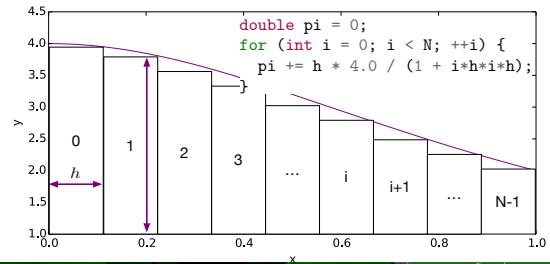
NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
GRID COMPUTINGUNIVERSITY OF
WASHINGTON**Numerical Quadrature**

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
GRID COMPUTINGUNIVERSITY OF
WASHINGTON**Numerical Quadrature**

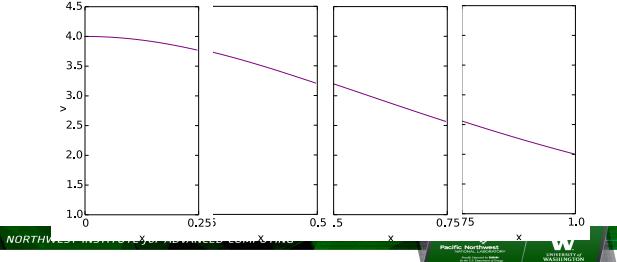
NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
GRID COMPUTINGUNIVERSITY OF
WASHINGTON**Numerical Quadrature (Sequential)**

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
GRID COMPUTINGUNIVERSITY OF
WASHINGTON**Finding Concurrency**

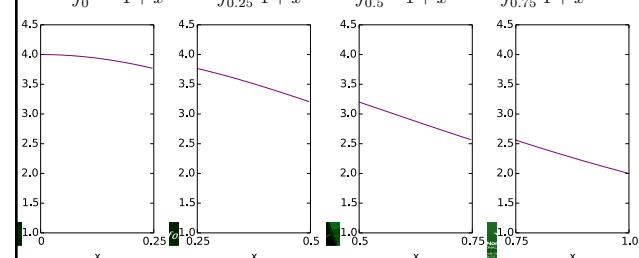
$$\pi = \int_0^{0.25} \frac{4}{1+x^2} dx + \int_{0.25}^{0.5} \frac{4}{1+x^2} dx + \int_{0.5}^{0.75} \frac{4}{1+x^2} dx + \int_{0.75}^1 \frac{4}{1+x^2} dx$$

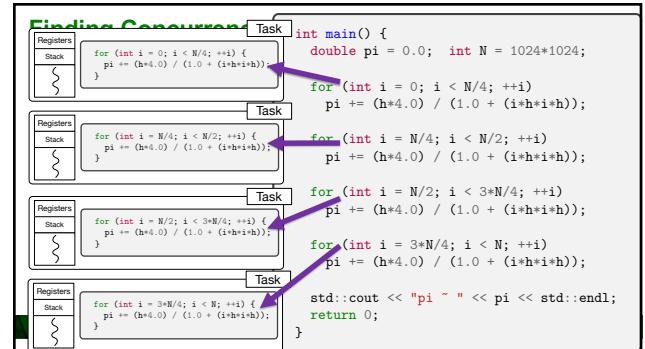
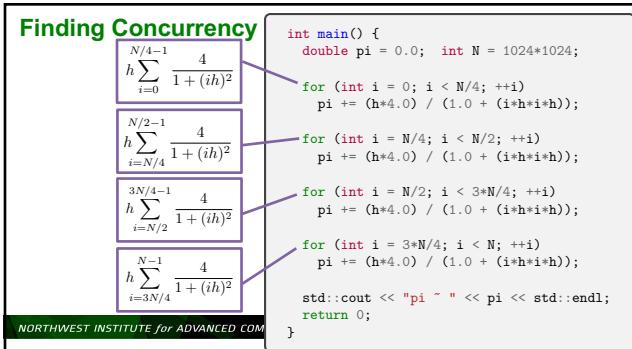
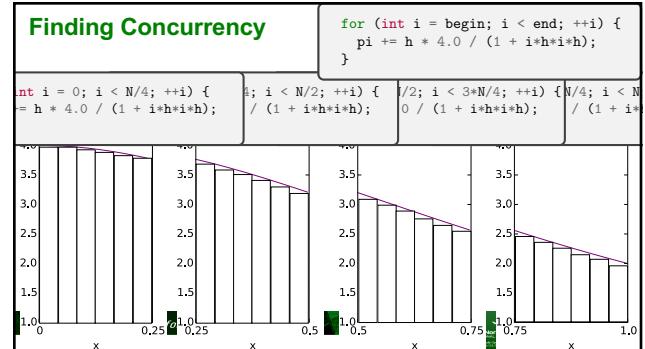
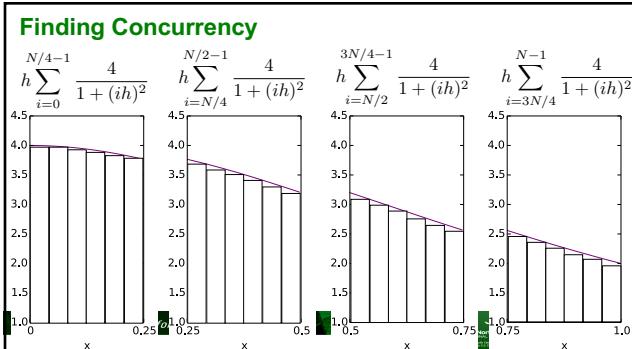
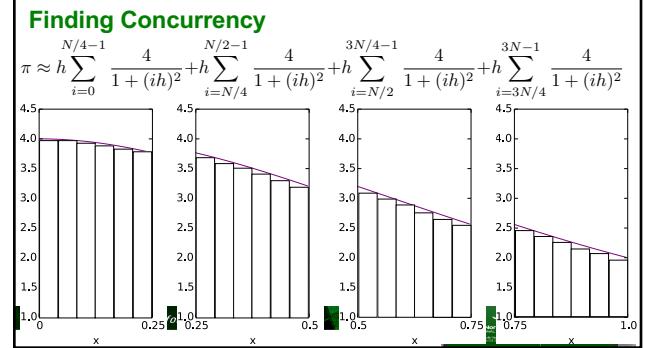
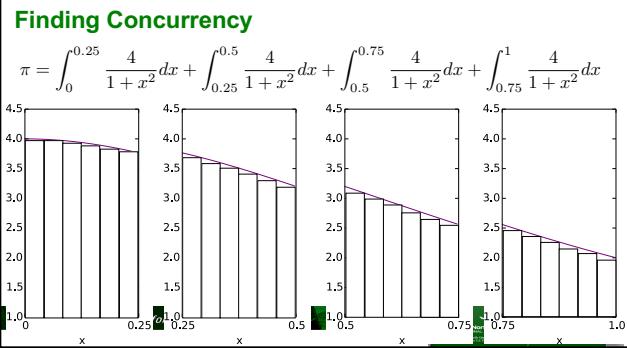


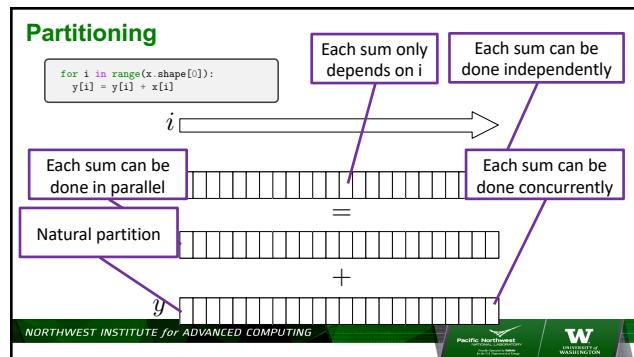
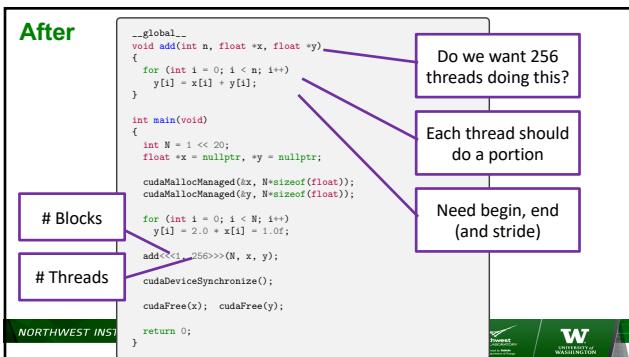
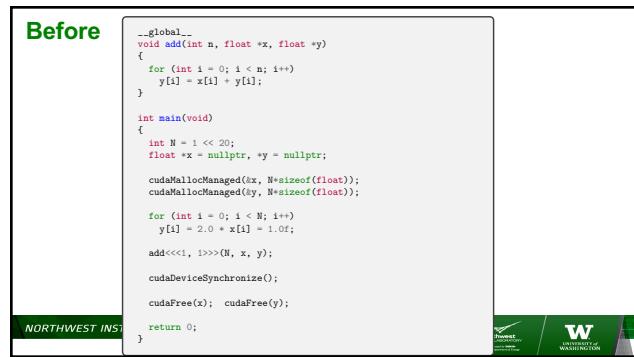
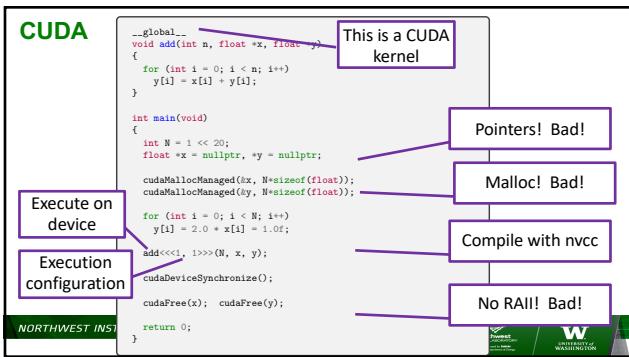
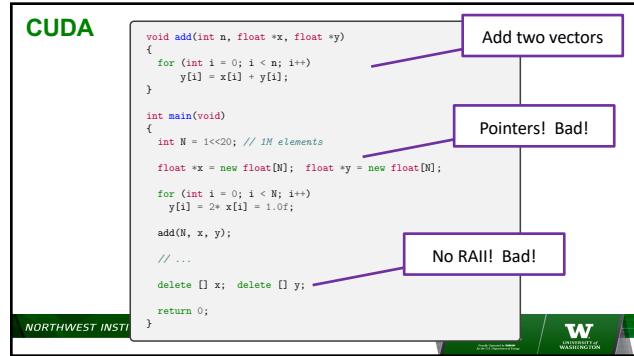
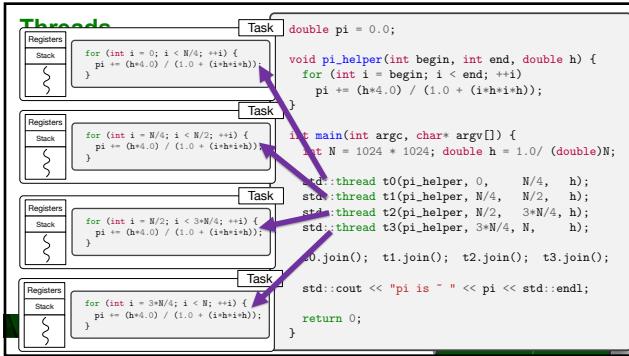
NORTHWEST INSTITUTE for ADVANCED COMPUTING

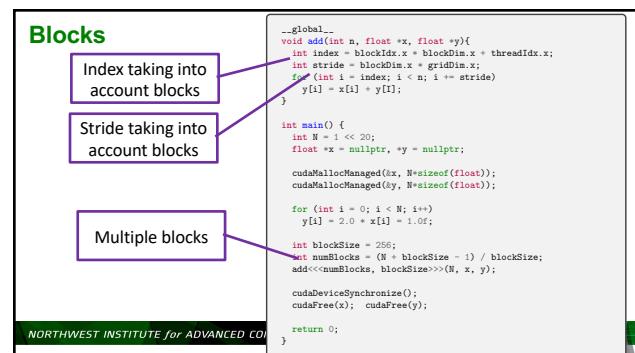
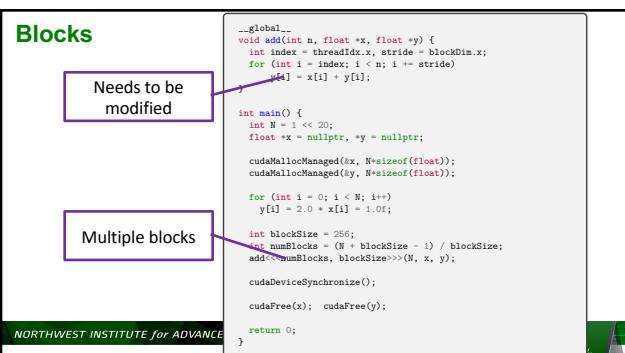
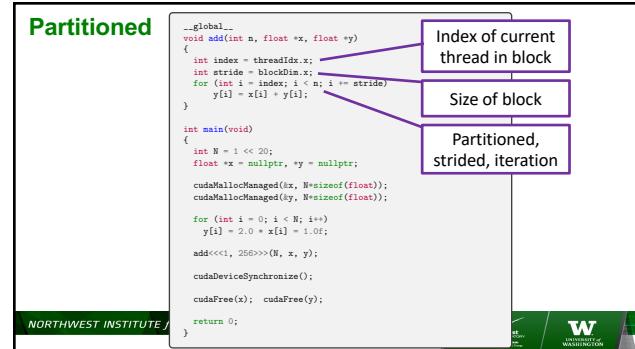
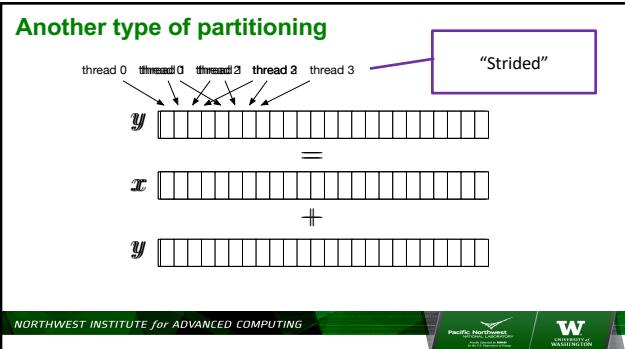
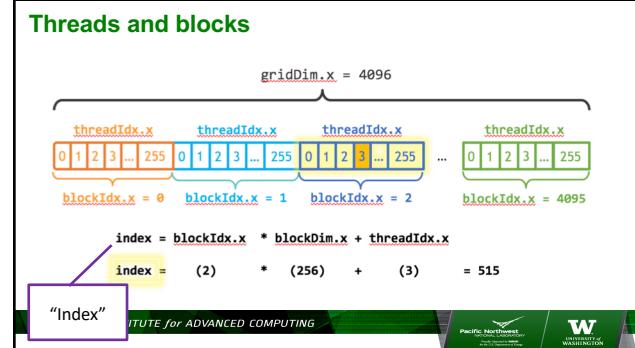
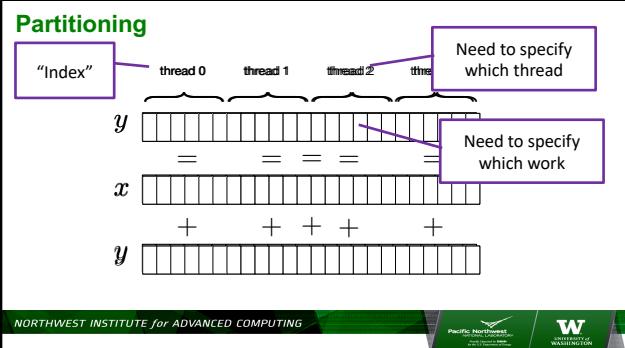
Pacific Northwest
GRID COMPUTINGUNIVERSITY OF
WASHINGTON**Finding Concurrency**

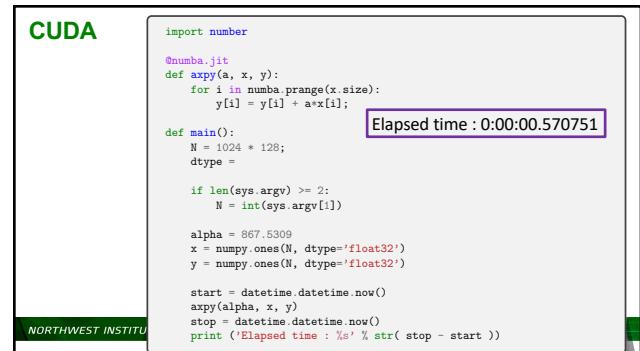
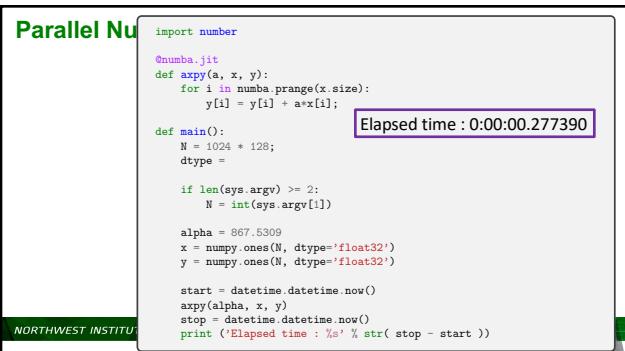
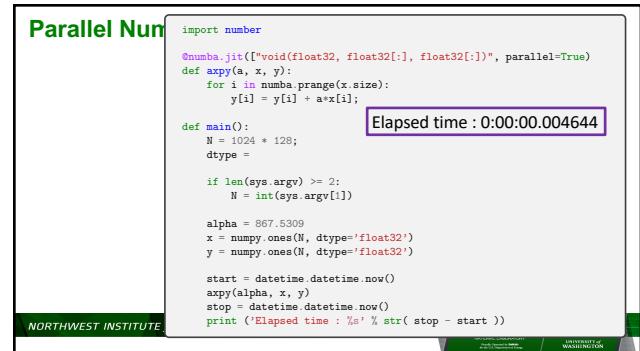
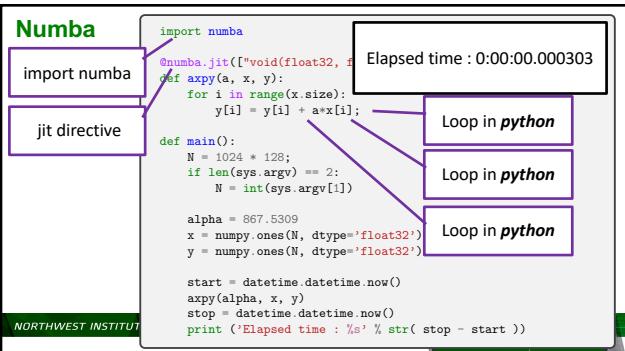
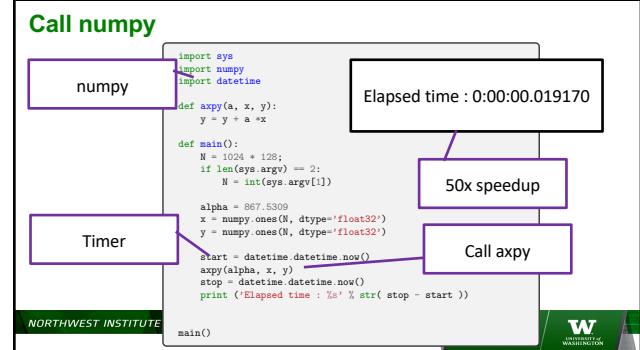
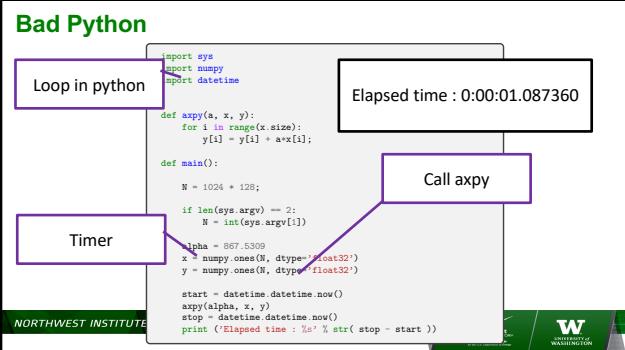
$$\pi = \int_0^{0.25} \frac{4}{1+x^2} dx + \int_{0.25}^{0.5} \frac{4}{1+x^2} dx + \int_{0.5}^{0.75} \frac{4}{1+x^2} dx + \int_{0.75}^1 \frac{4}{1+x^2} dx$$











CUDA

```

import numba
from numba import cuda

@cuda.jit
def axpy(x, y):
    tx = cuda.threadIdx.x
    ty = cuda.blockIdx.x
    bw = cuda.blockDim.x
    pos = tx * ty * N
    if pos < x.size:
        y[pos] = y[pos] + a*x[pos]

def main():
    N = 1024 * 128;

    if len(sys.argv) >= 2:
        N = int(sys.argv[1])

    alpha = 1.00000001;
    x = numpy.ones(N, dtype=dtype)
    y = numpy.ones(N, dtype=dtype)

    y[0] = 0; y[N-1] = 0;
    start = datetime.datetime.now()
    axpy((N-1)/250, alpha, x, y)
    stop = datetime.datetime.now()
    print('Elapsed time : %s' % (str( stop - start )))

```

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Performance

	V.1	V.2	Laptop (GeForce GT 750M)		Server (Tesla K80)	
Version			Time	Bandwidth	Time	Bandwidth
1 CUDA Thread			411ms	30.6 MB/s	463ms	27.2 MB/s
1 CUDA Block			3.2ms	3.9 GB/s	2.7ms	4.7 GB/s
Many CUDA Blocks			0.68ms	18.5 GB/s	0.094ms	134 GB/s
	V.3				5000X	!!

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest National Laboratory
DOE Office of Science
University of Washington

W
UNIVERSITY OF WASHINGTON

Summary (so far)

- Transfer data from CPU memory to GPU memory
- Compute in GPU
- Transfer data back
- Computation is done with huge number of identical processing units
- Each PU executes exactly same kernel
- Partition the data among the processors – use thread idx, block idx, and block dimension to compute a unique (global)
- Use global index in the kernel to determine data each thread should work on

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
National Laboratory

W
UNIVERSITY OF WASHINGTON

Other GPU programming systems

- OpenACC
- OpenCL
- SyCL
- CuSP
- CuBLAS
- Halide
- Etc.

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
National Laboratory

W
UNIVERSITY OF WASHINGTON

Thank you!

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
National Laboratory

W
UNIVERSITY OF WASHINGTON

Creative Commons BY-NC-SA 4.0 License



© Andrew Lumsdaine, 2017-2018

Except where otherwise noted, this work is licensed under

<https://creativecommons.org/licenses/by-nc-sa/4.0/>

Cuda and Thrust programming examples © Nvidia

NORTHWEST INSTITUTE for ADVANCED COMPUTING

Pacific Northwest
National Laboratory

W
UNIVERSITY OF WASHINGTON

