

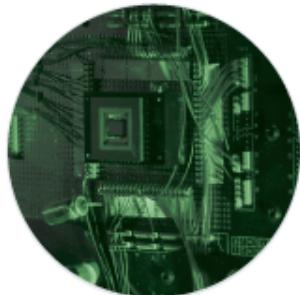


A. JAMES CLARK
SCHOOL OF ENGINEERING

A Standalone Package for Bringing Graphics Processor Acceleration to GNU Radio: *GRGPU*

William Plishker

University of Maryland
plishker@umd.edu





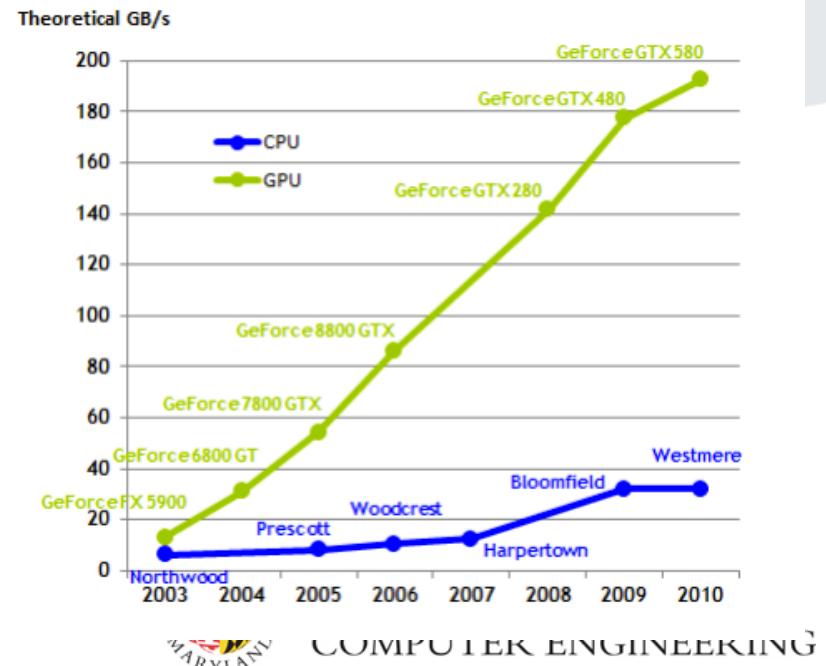
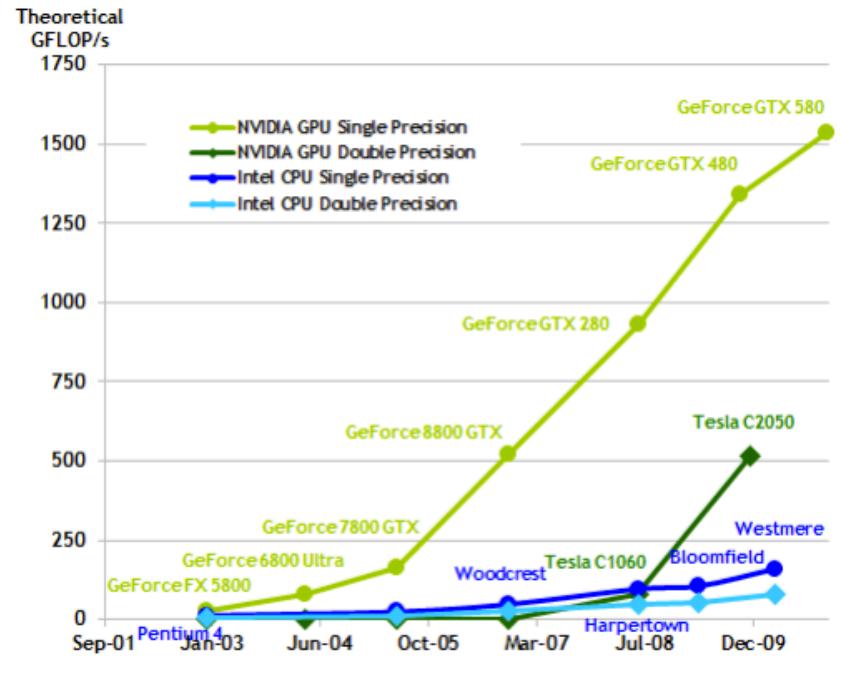
Outline

- Introduction
- GPU Background
- Graphics Processor Integration with GNU Radio
- Evaluation
- Summary



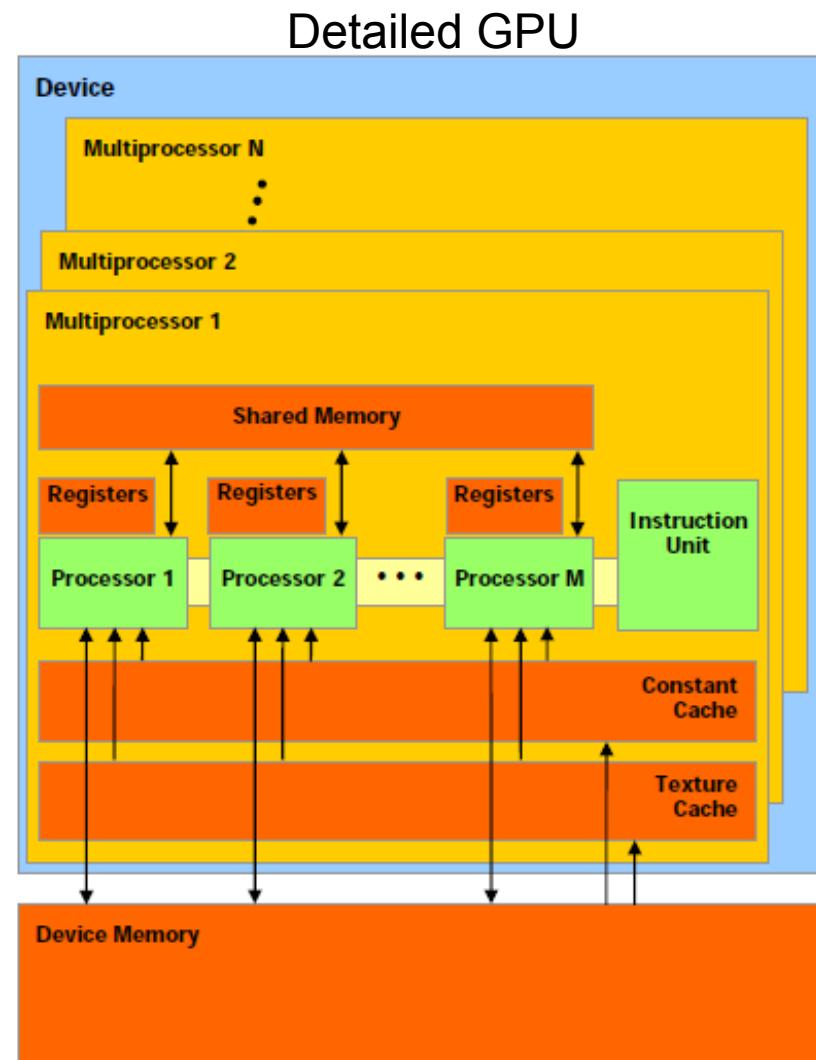
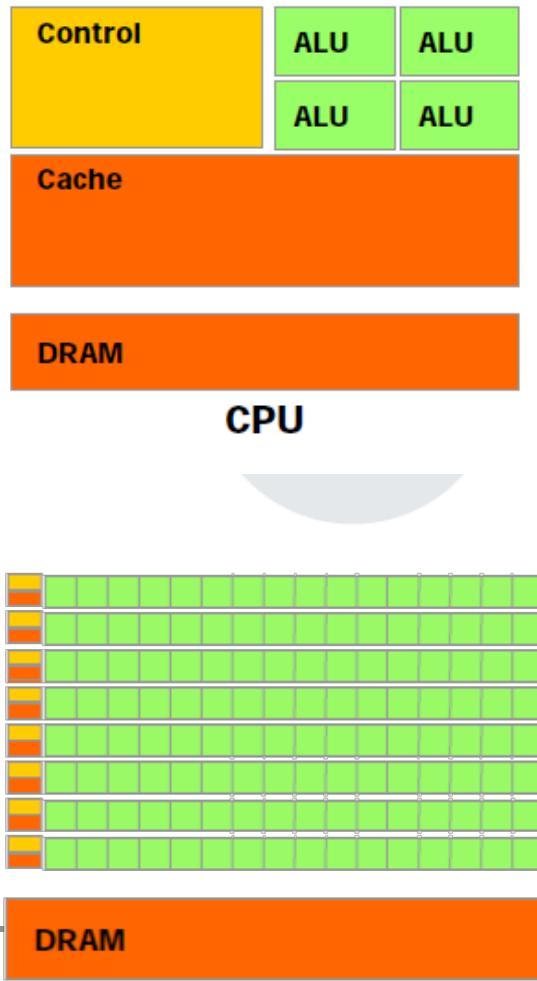
Motivation

- High-Performance prototyping with GPUs is tempting...
 - *Potentially* fast design times
 - Flop/\$
- ...but hard.
 - Learn a new language
 - Expose or refactor parallelism





The GPU Difference





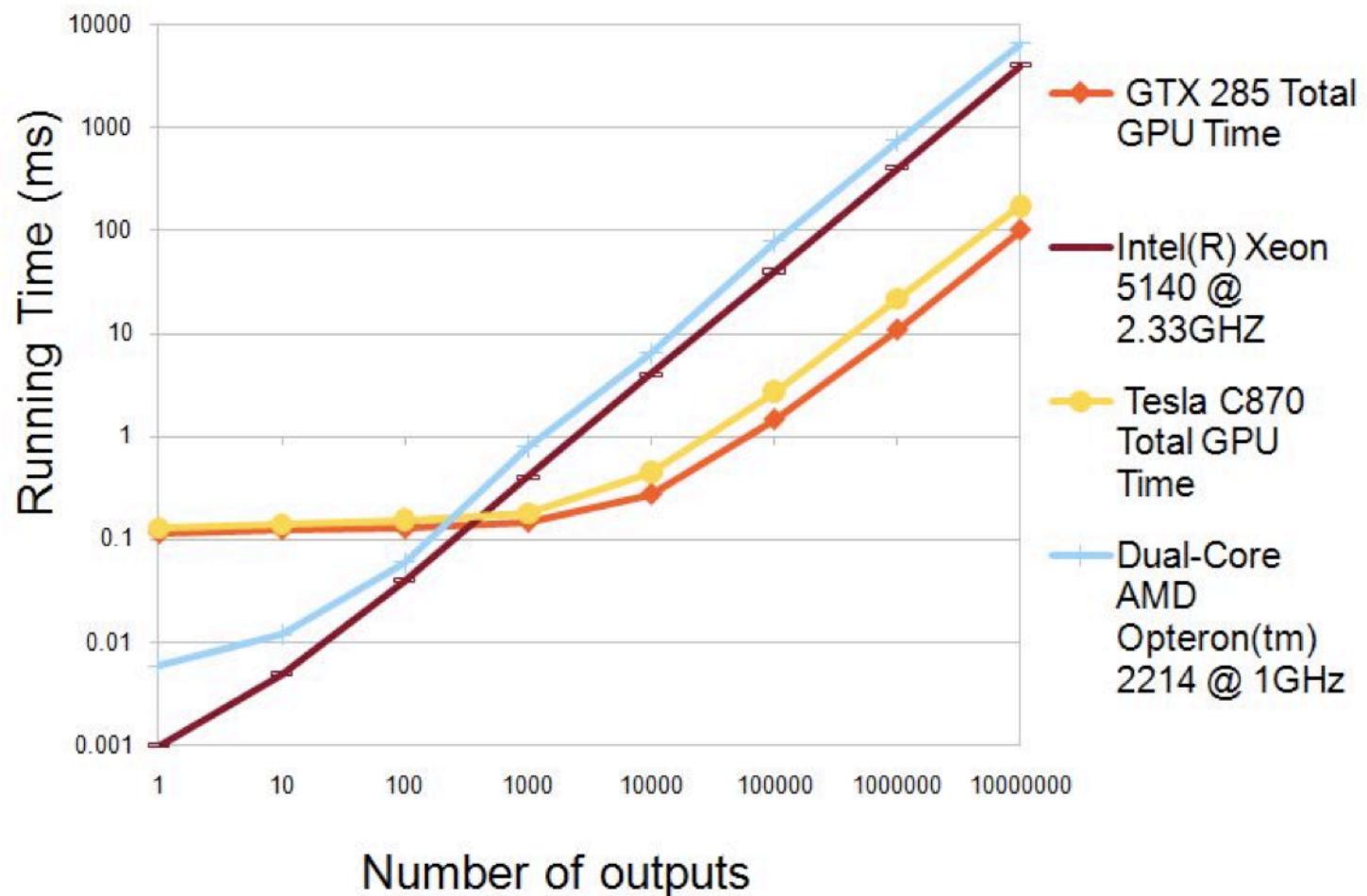
A typical high end modern GPU: GTX 580

- Released 10/2010
- 3B transistors, 40nm
- 512 Cores @700MHz
- 1.5 TFLOPs
- 192 GB/s Memory BW
- 244W
- \$500





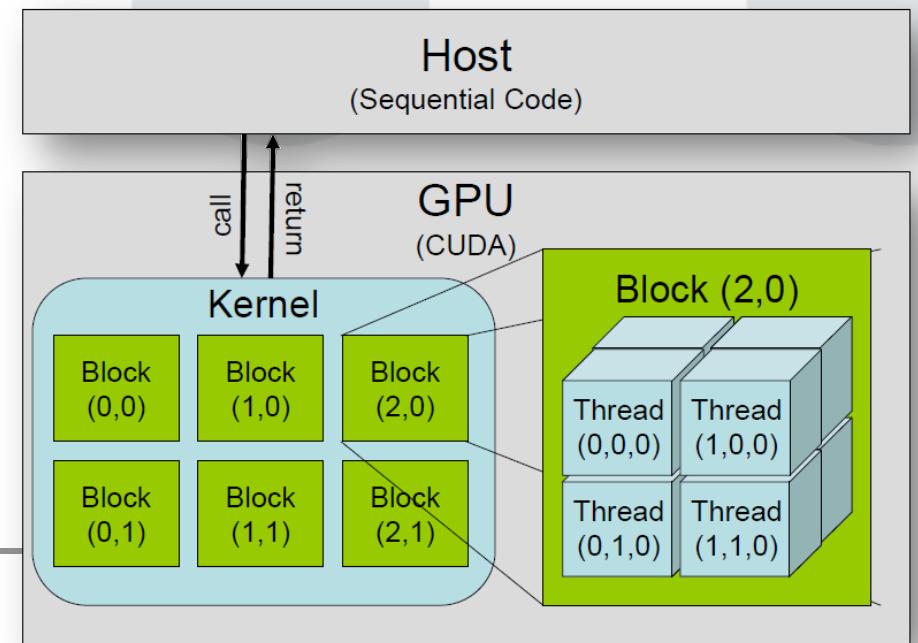
CUDA FIR vs. Stock FIR





A brief history of GPU Programming

- OpenGL (1992)/DirectX(1995) + shader programming in assembly
- Cg, GLSL, HLSL (early 2000s)
- Sh, Brook (2004)
- CUDA (2006)
- OpenCL (2008)





Related Work – MATLAB

A screenshot of the MATLAB 7.11.0 (R2010b) interface. The workspace browser on the left shows variables A, b, f, f_gpu, x, and x_gpu. The command window on the right displays MATLAB code demonstrating GPU operations. The code creates GPU arrays, uses GPU-enabled functions like \ and fft, and gathers data back into the MATLAB workspace.

```
>>
>> % Create arrays that reside on the GPU
>> A = gpuArray(rand(1000, 1000));
>> b = gpuArray(rand(1000, 1));
>>
>> % Use GPU-enabled MATLAB functions
>> x_gpu = A \ b; % "\" is GPU-enabled
>>
>> f_gpu = fft(A);
>>
>> % Bring data back from GPU memory into MATLAB workspace
>> x = gather(x_gpu);
>> f = gather(f_gpu);
fx>> |
```



Related Work – PyCUDA

```
import pycuda.autoinit
import pycuda.driver as drv
import numpy

from pycuda.compiler import SourceModule
mod = SourceModule("""
__global__ void multiply_them(float *dest, float *a, float *b)
{
    const int i = threadIdx.x;
    dest[i] = a[i] * b[i];
}
""")

multiply_them = mod.get_function("multiply_them")

a = numpy.random.randn(400).astype(numpy.float32)
b = numpy.random.randn(400).astype(numpy.float32)

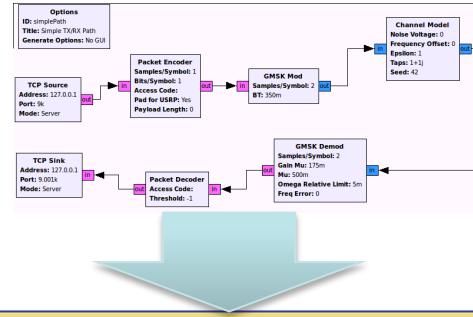
dest = numpy.zeros_like(a)
multiply_them(
            drv.Out(dest), drv.In(a), drv.In(b),
            block=(400,1,1), grid=(1,1))

print dest-a*b
```



Background – GNU Radio Design Flow

GNU Radio Companion (GRC)



GNU Radio Python Flowgraph

```
class top(gr.top_block):
    def __init__(self):
        gr.top_block.__init__(self)
        ntaps = 256

        # Something vaguely like floating point ops
        self.flop = 2 * ntaps * options.npipelines

        src = gr.null_source(gr.sizeof_float)
        head = gr.head(gr.sizeof_float, int(options.nsamples))
        self.connect(src, head)

        for n in range(options.npipelines):
            self.connect(head, pipeline(options.nstages, ntaps))
```



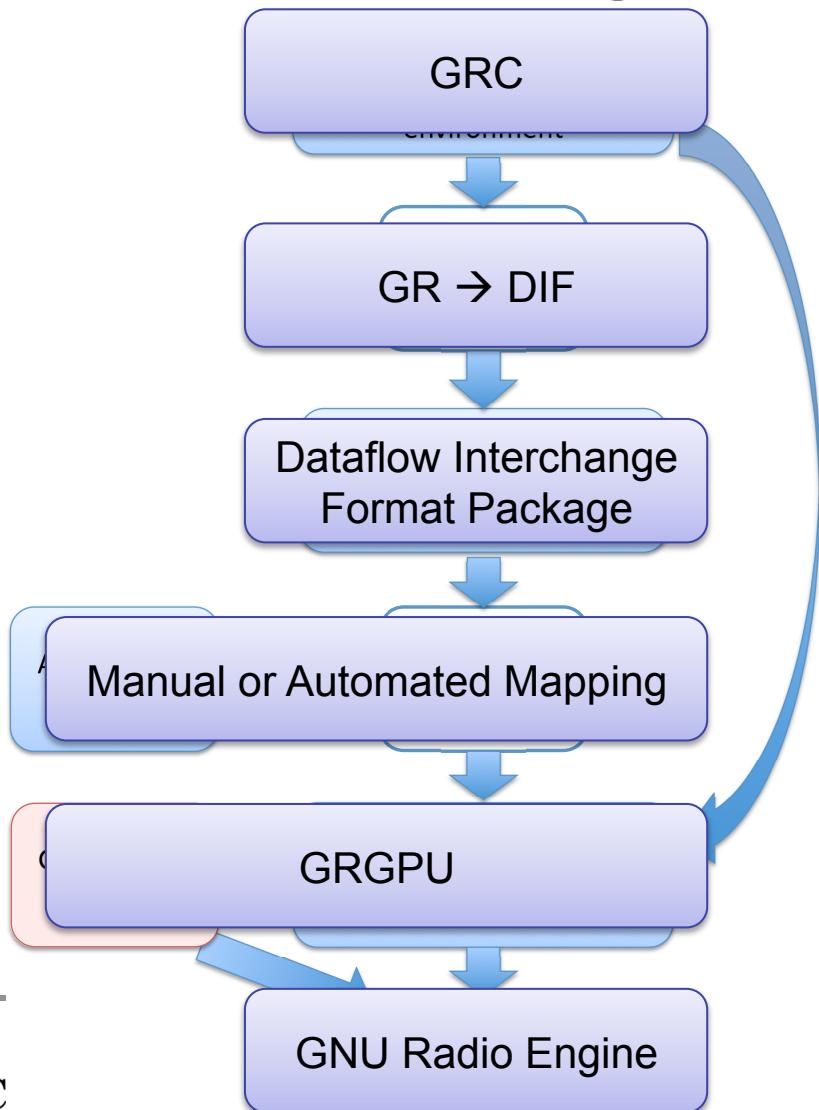
*Universal Software
Radio Peripheral
(USRP)*



DEPARTMENT OF
ELECTRICAL &
COMPUTER ENGINEERING



Proposed Design Flow

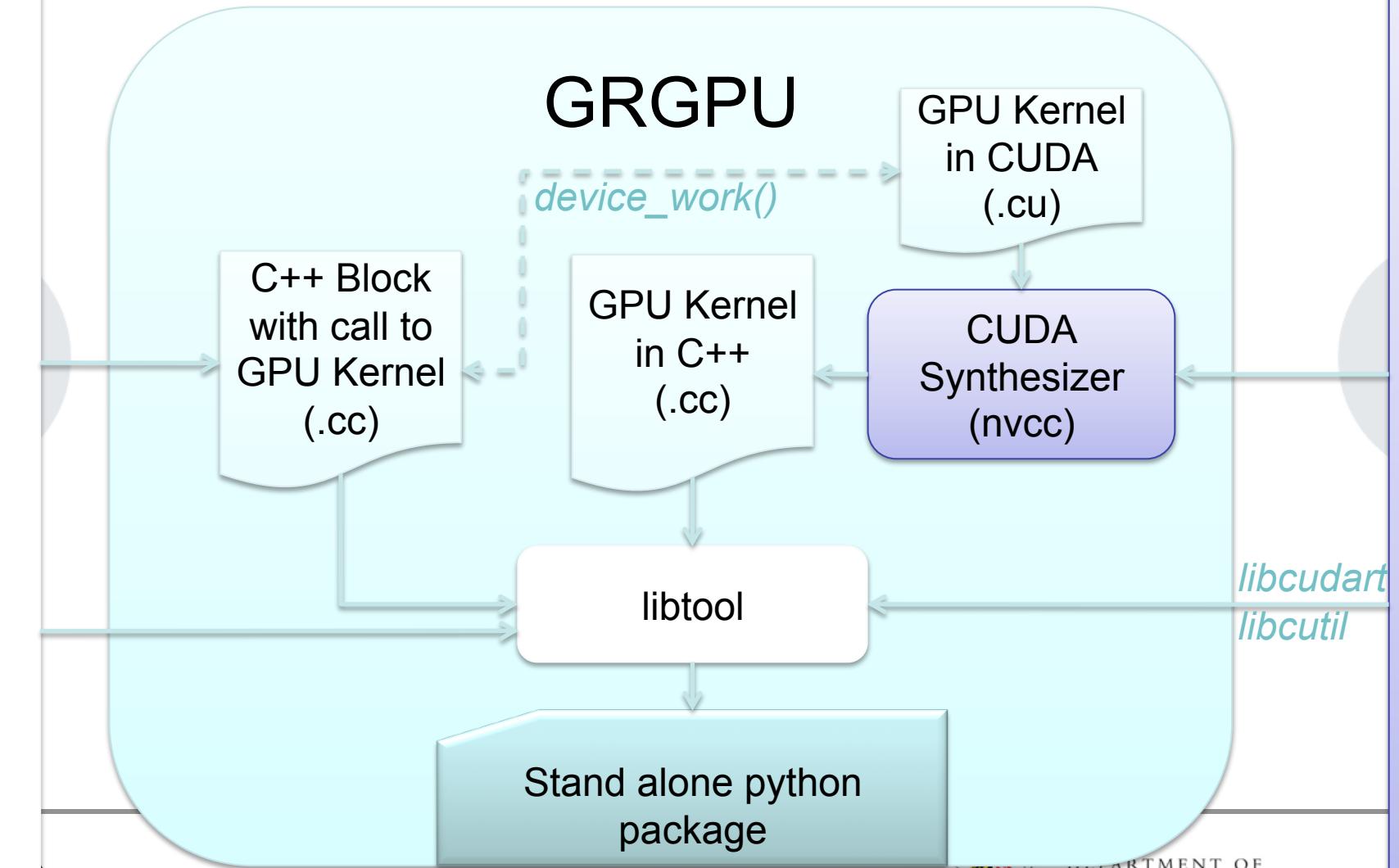


- Start from an unambiguous application description based on the application area itself
 - Dataflow
- Use tools to optimize scheduling, assignment
- Generate an accelerated functional description



GNU Radio for Graphics Processing Units (GRGPU)

GNU Radio





Build Process

- make-cuda : .cu → .cc
 - nvcc as a software synthesizer
 - transform cuda code into target specific GPU code
 - wrap functions in extern “C”
- make
 - Resulting cc's are just more inputs to existing build structure
 - link against CUDA runtime (cudart) and CUDA utils (cutil)



Using GRGPU

- Used just like gr in GNU Radio flowgraphs
 - import grgpu
 - op = grgpu.op_cuda()
- Dependent blocks
 - Host/device data movement handled by other blocks
 - Has device pointers as interface



GRGPU Example

```
class GPU_Thread(threading.Thread):
    """GPU Thread"""
    def init(self):
        self.tb = gr.top_block()
        src_data = list(math.sin(x) for x in range(TOTAL_LENGTH))
        src = gr.vector_source_f(src_data)

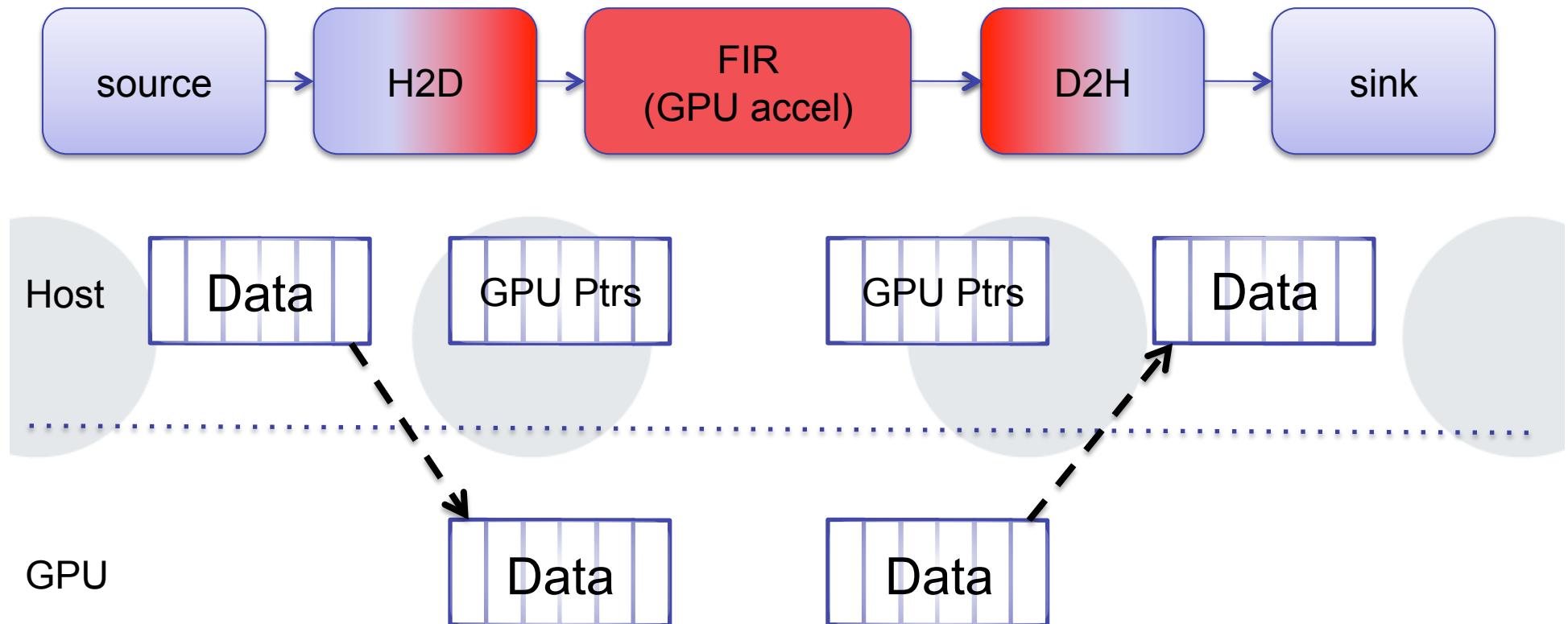
        h2d = grgpu.h2d_cuda()
        taps = [0.1]*60
        op = grgpu.fir_filter_fff_cuda(taps)
        d2h = grgpu.d2h_cuda()

        self.dst = gr.vector_sink_f()
        self.tb.connect(src, h2d)
        self.tb.connect(h2d, op)
        self.tb.connect(op, d2h)
        self.tb.connect(d2h, self.dst)
```





GRGPU Example





Other GRGPU features

- Fifo reads/writes are chunky and have history
 - Designed to make use of burst transfers
 - Creates in data contiguous memory for kernels
- Growing library
 - add_const, mul_const, FIR, resampler, FFT
 - floats, complexes
- Typical GR conveniences
 - Create new GRGPU block script
 - Doxygen, unit tests, GR build process

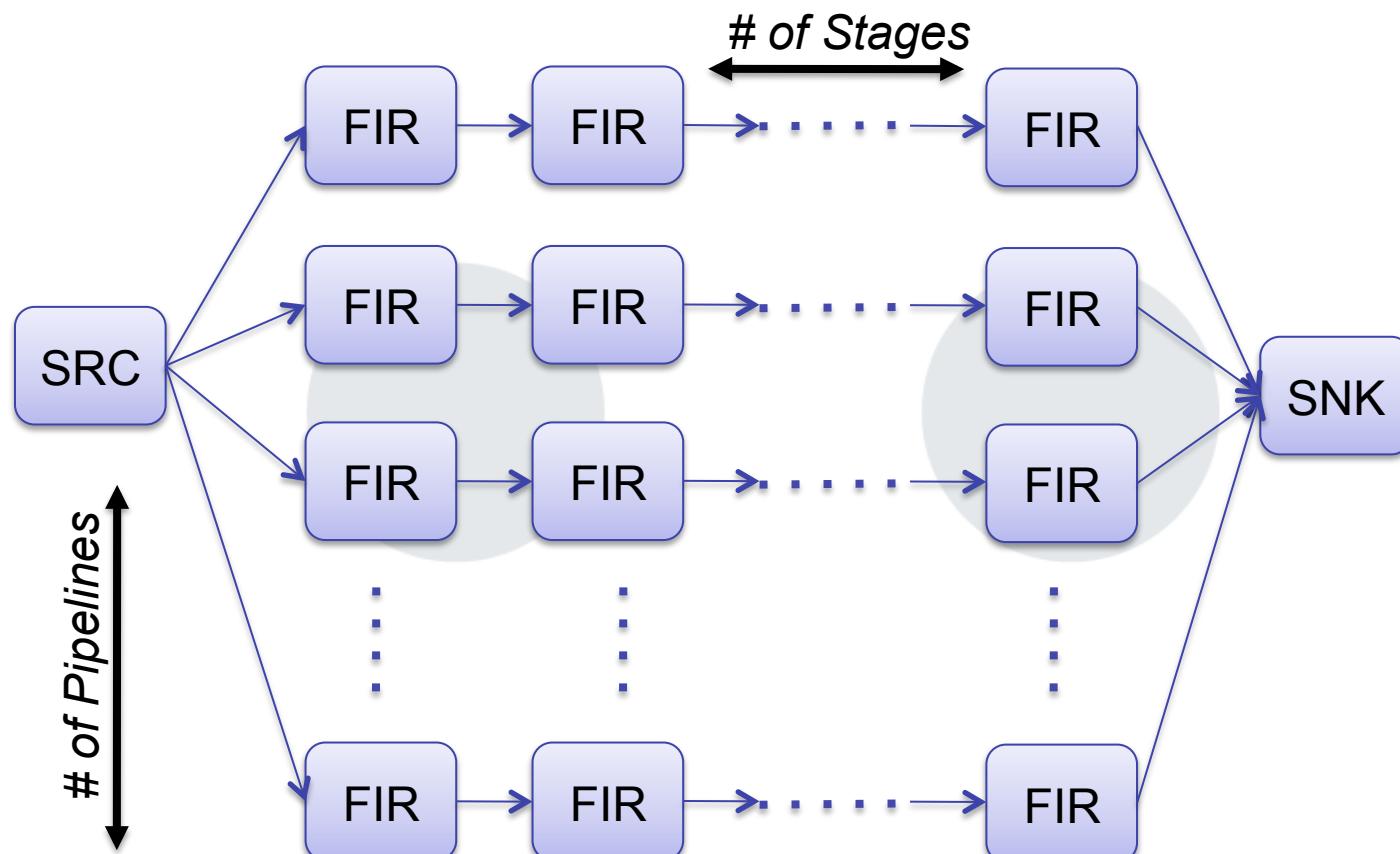


Evaluation

- Benchmark GPU options with GNU Radio multicore benchmark (mp-sched)
 - 2 Intel Xeons (3GHz) and 1 GTX 260
- GPU and GPP accelerated FIR
 - CUDA and SSE acceleration
- Vary: size and assignment
- Measure: total execution time

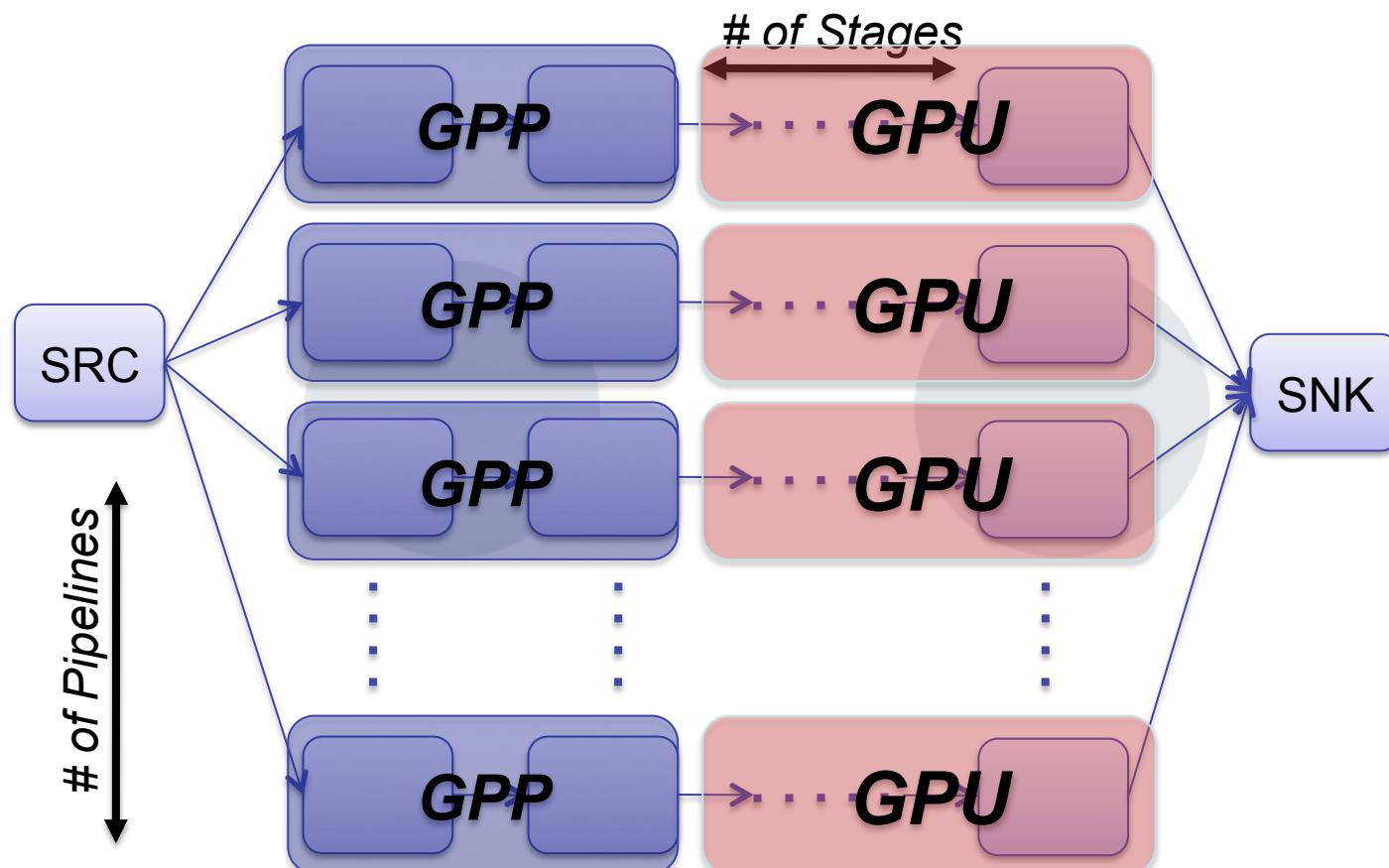


Evaluation Application



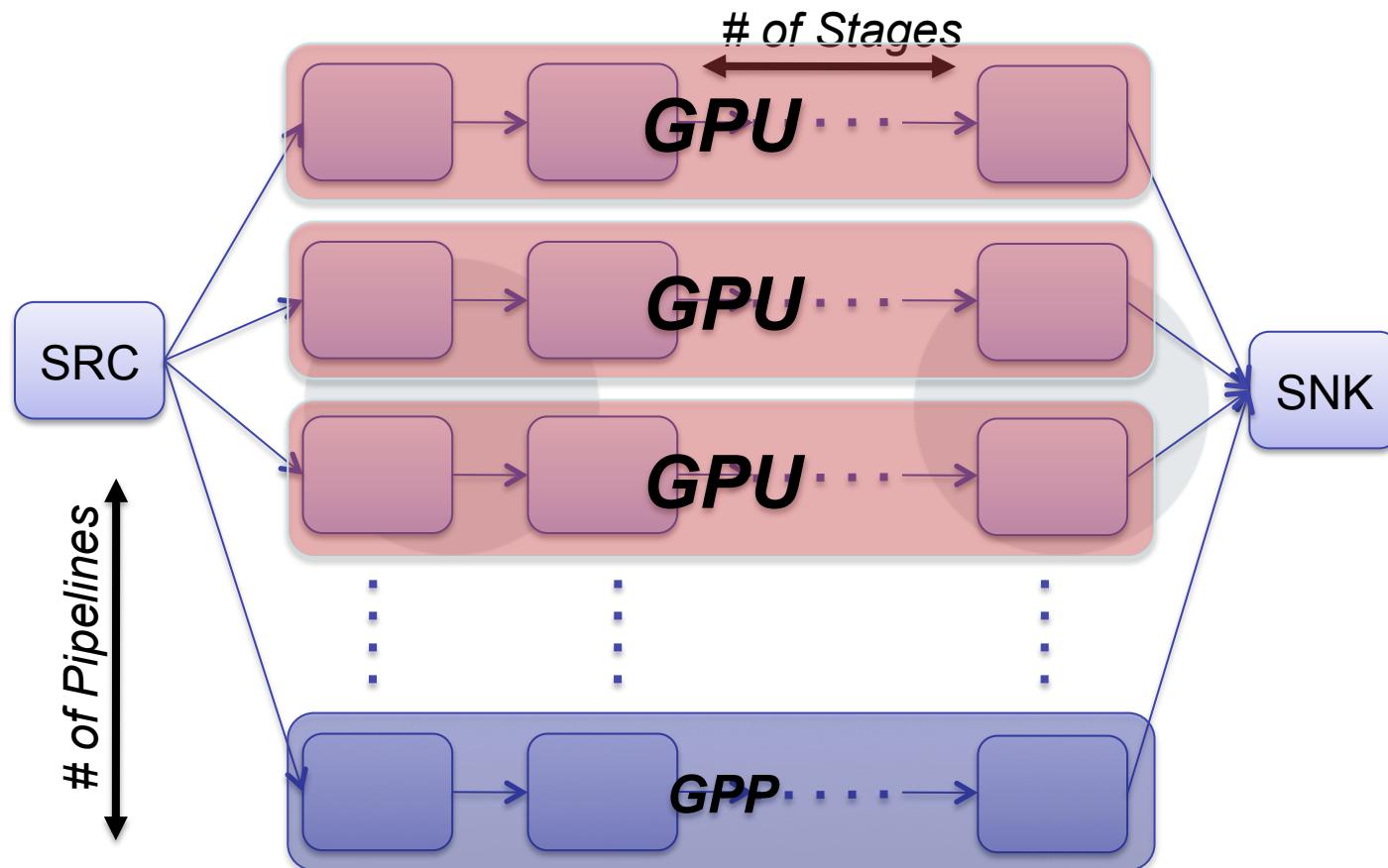


Evaluation Application – “Half CPU”



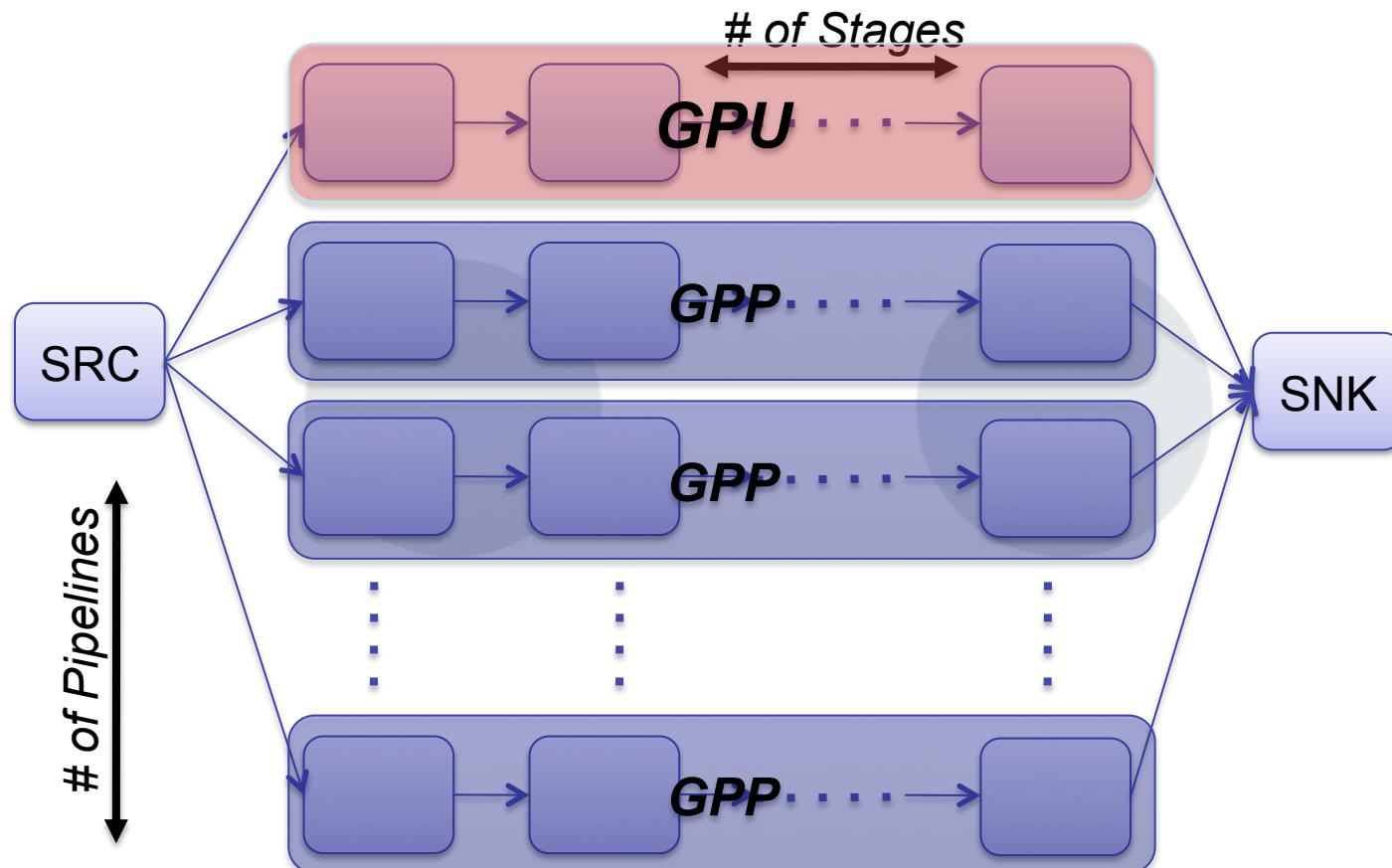


Evaluation Application – “One GPP”



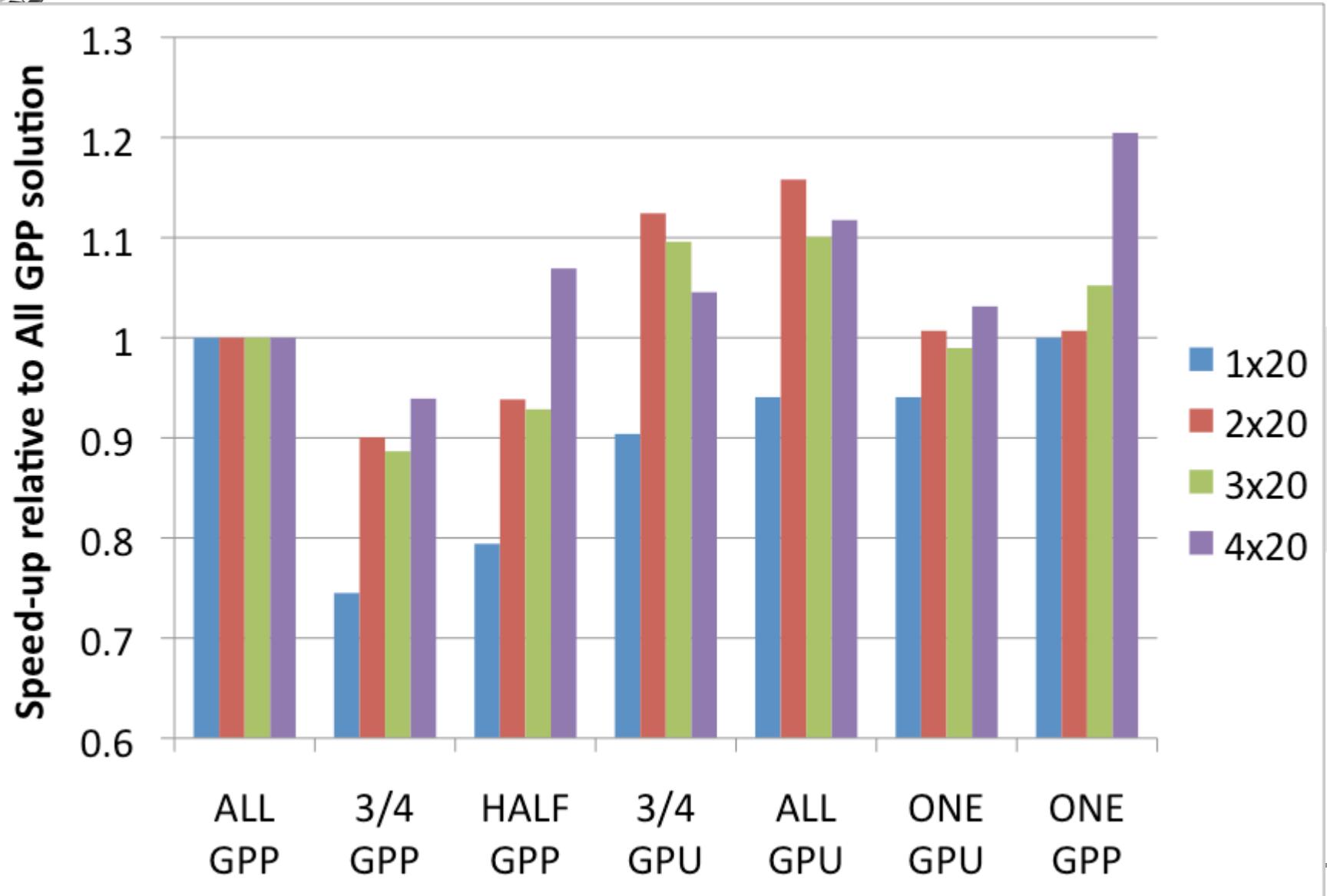


Evaluation Application – “One GPU”





Results





Conclusions

- GRGPU provides GPU acceleration to GNU Radio as a stand alone library
 - Facilitates design space exploration
- Status
 - Planned to be contributed back
 - Few but growing number of CUDA accelerated blocks
- Ongoing Work
 - Automated Mapping
 - MultiGPU



Questions?

Citations:

W. Plishker, G. Zaki, S. S. Bhattacharyya, C. Clancy, and J. Kuykendall. Applying graphics processor acceleration in a software defined radio prototyping environment. In *Proceedings of the International Symposium on Rapid System Prototyping*, pages 67-72, Karlsruhe, Germany, May 2011.

G. Zaki, W. Plishker, T. OShea, N. McCarthy, C. Clancy, E. Blossom, and S. S. Bhattacharyya. Integration of dataflow optimization techniques into a software radio design framework. In *Proceedings of the IEEE Asilomar Conference on Signals, Systems, and Computers*, pages 243-247, Pacific Grove, California, November 2009.