



Accelerated Python

Mark Ebersole, NVIDIA
CUDA Educator





Once upon a time...

Past Massively Parallel Supercomputers



Thinking Machine



Goodyear MPP



MasPar



Cray 2

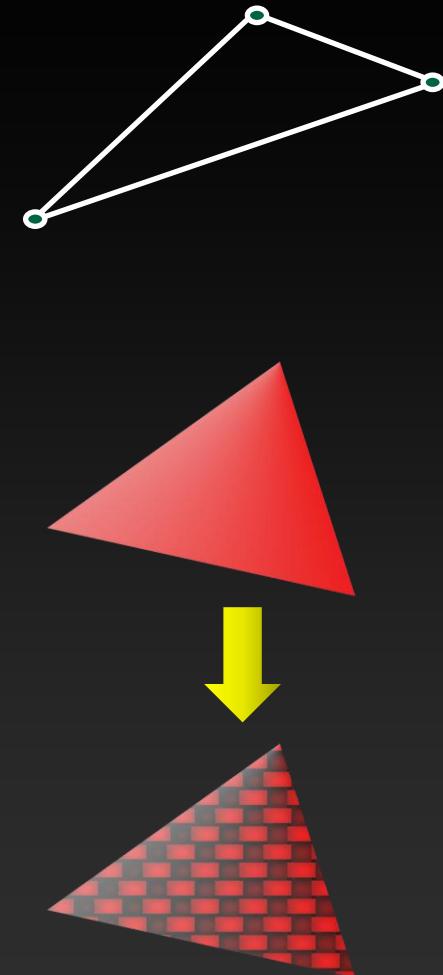
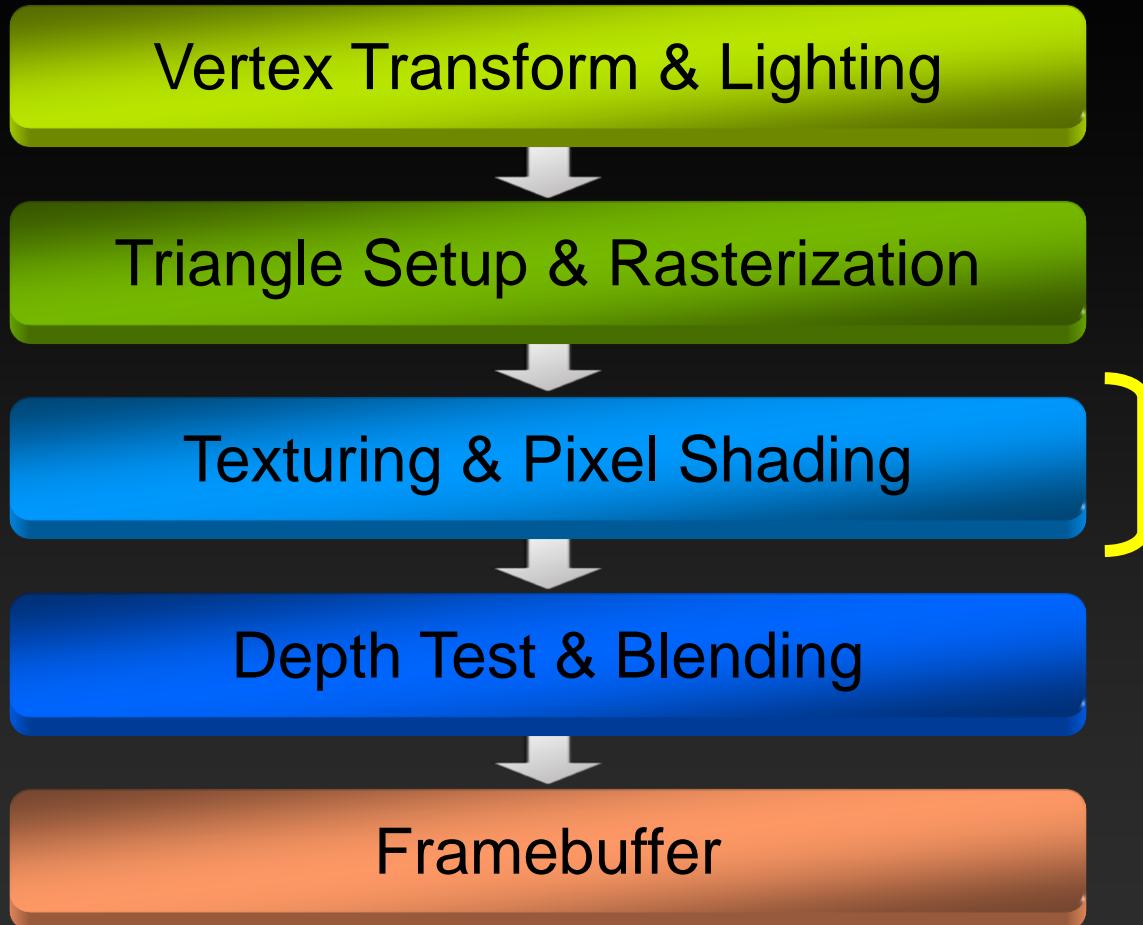


1.31 TFLOPS on
DGEMM

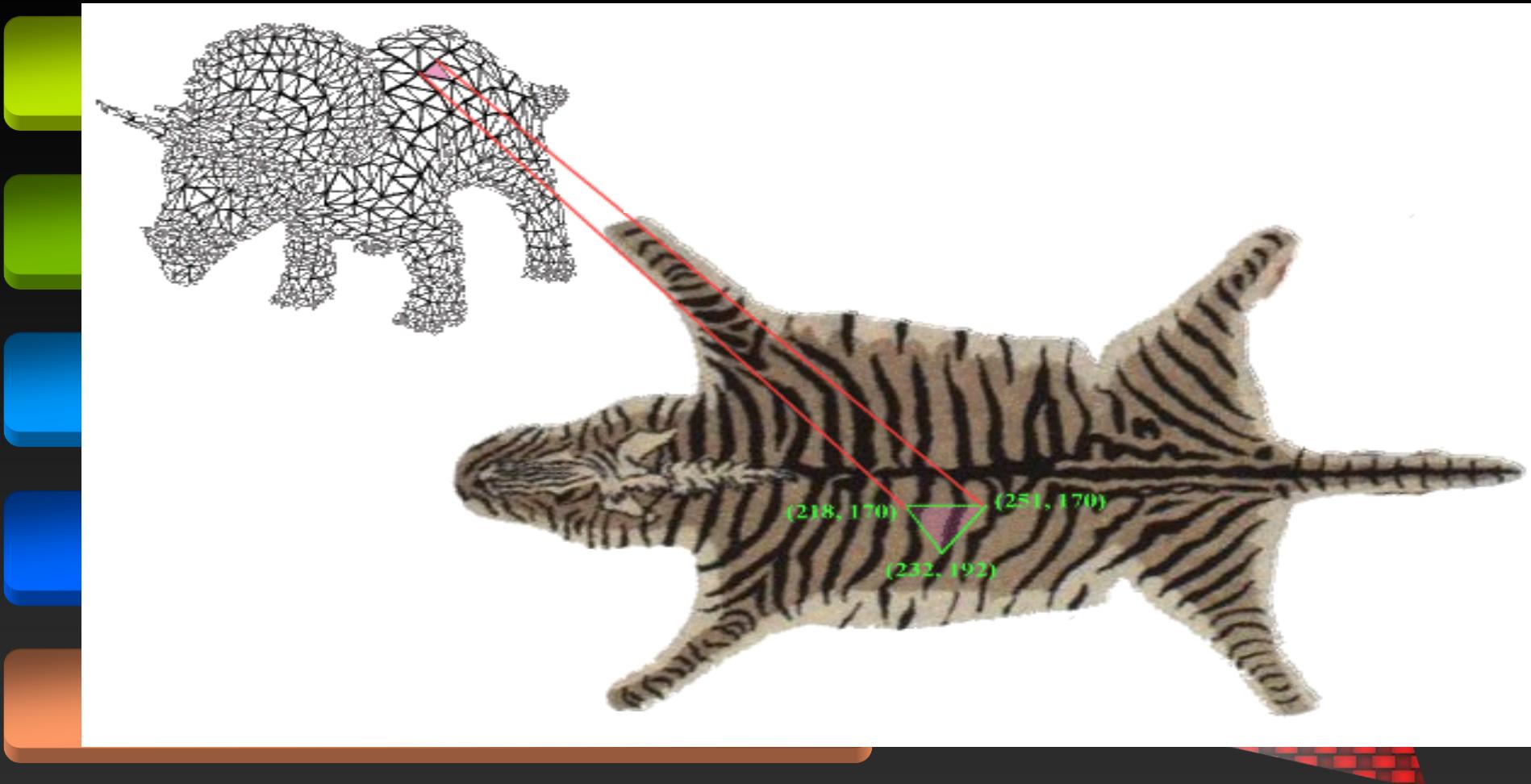




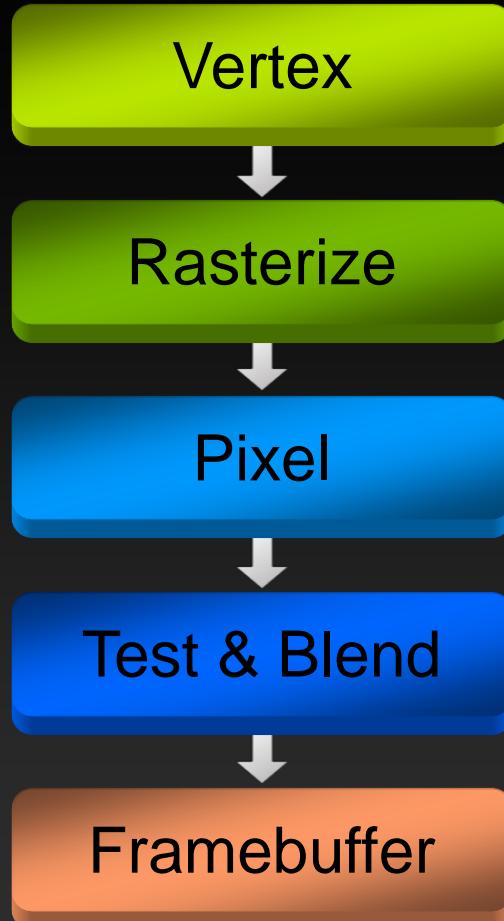
The Graphics Pipeline



The Graphics Pipeline

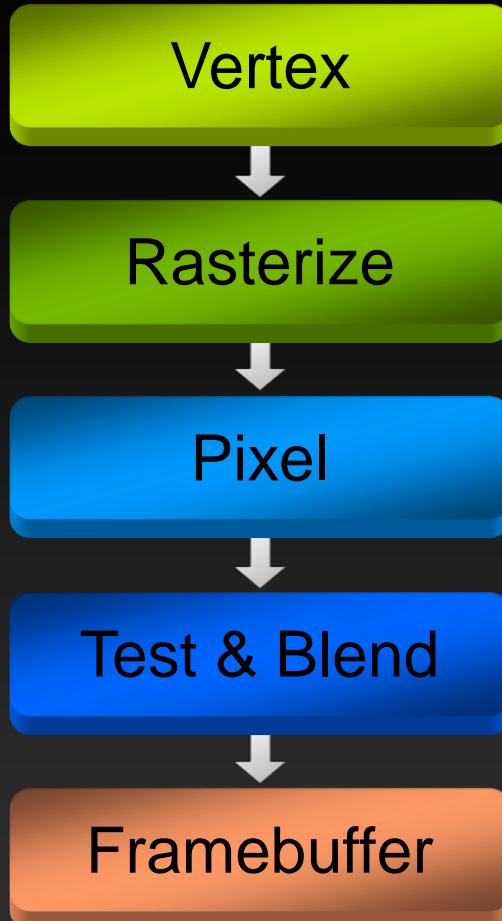


The Graphics Pipeline



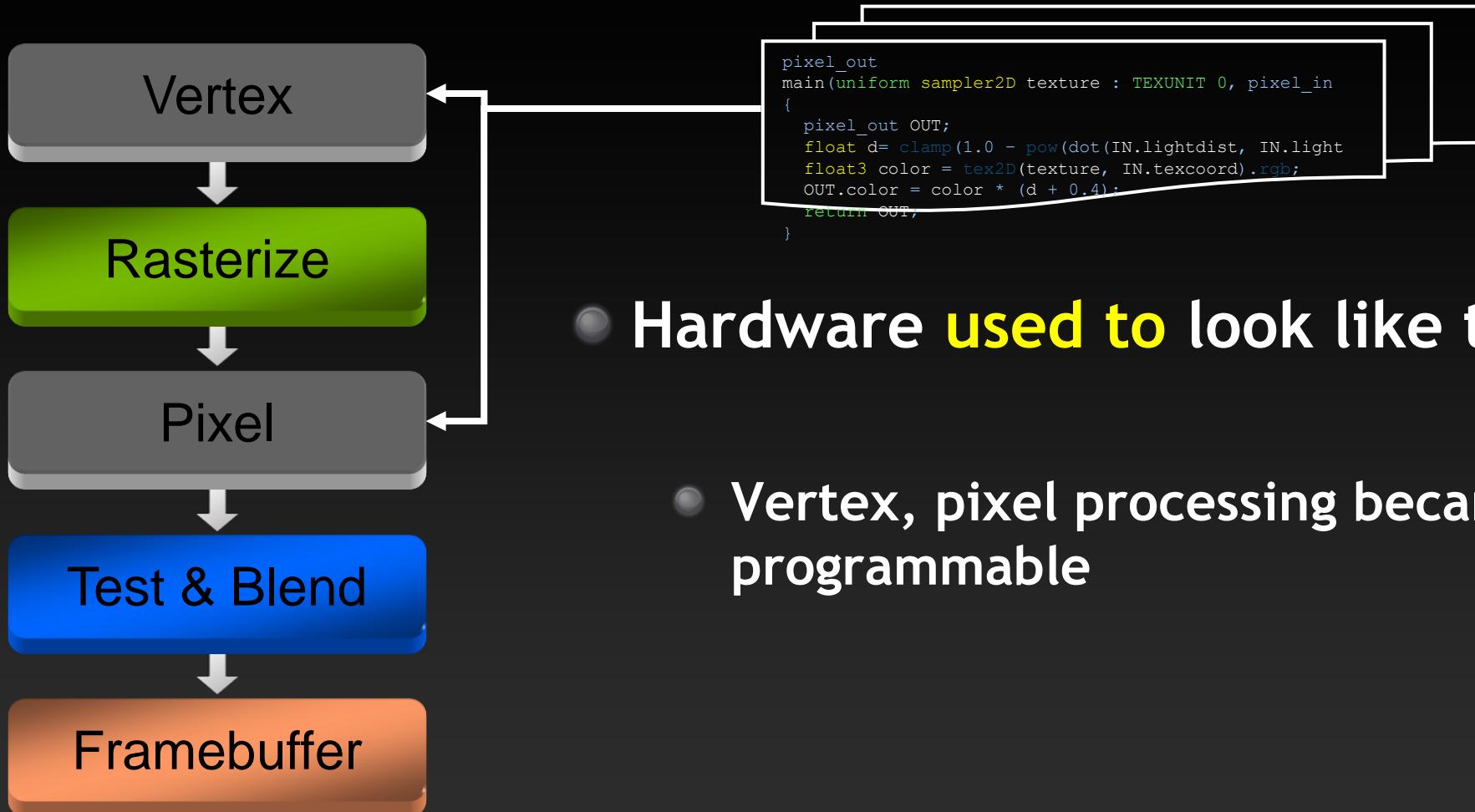
- Key abstraction of real-time graphics
- Hardware used to look like this
- One chip/board per stage
- Fixed data flow through pipeline

The Graphics Pipeline

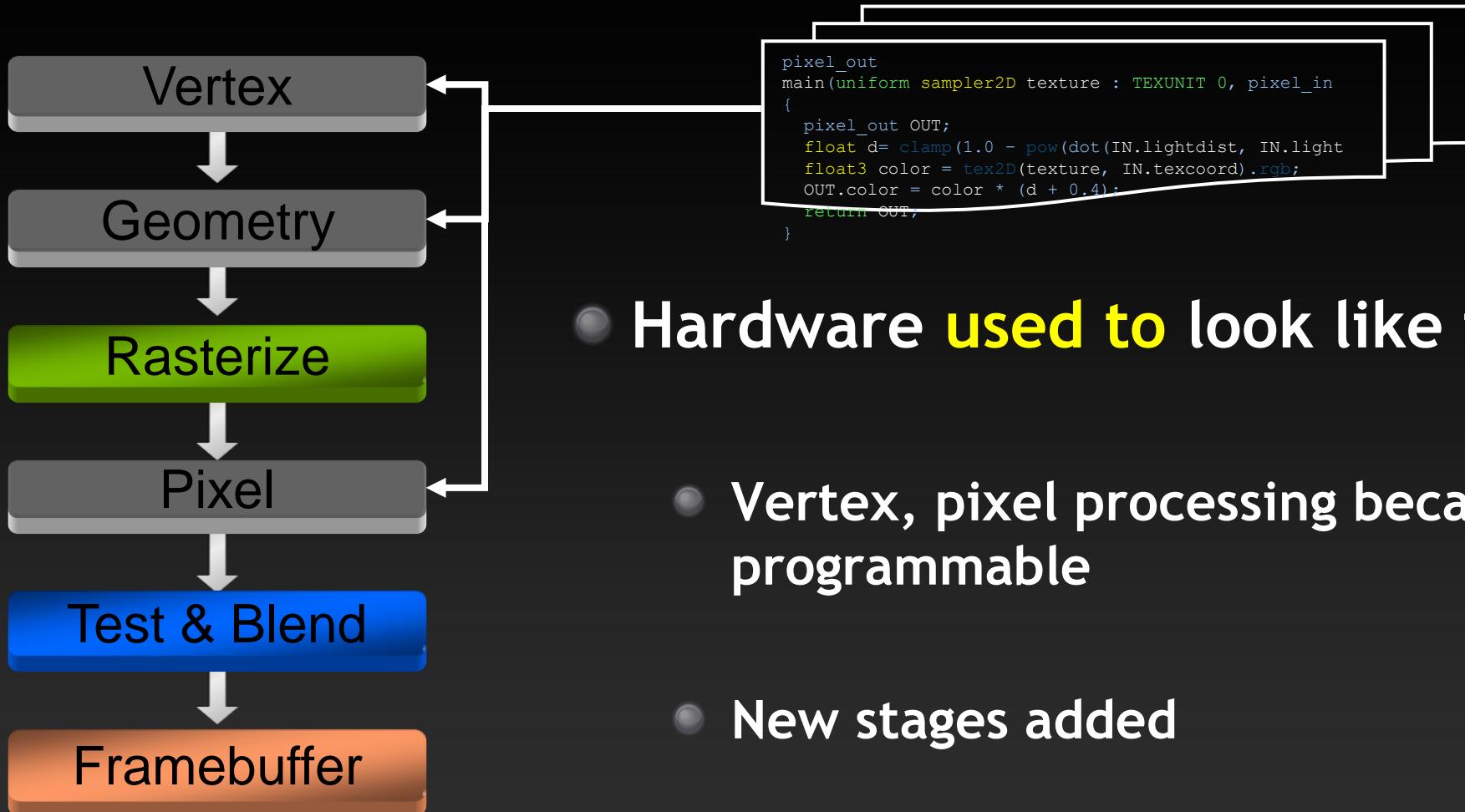


- Remains a useful abstraction
- Hardware used to look like this

The Graphics Pipeline

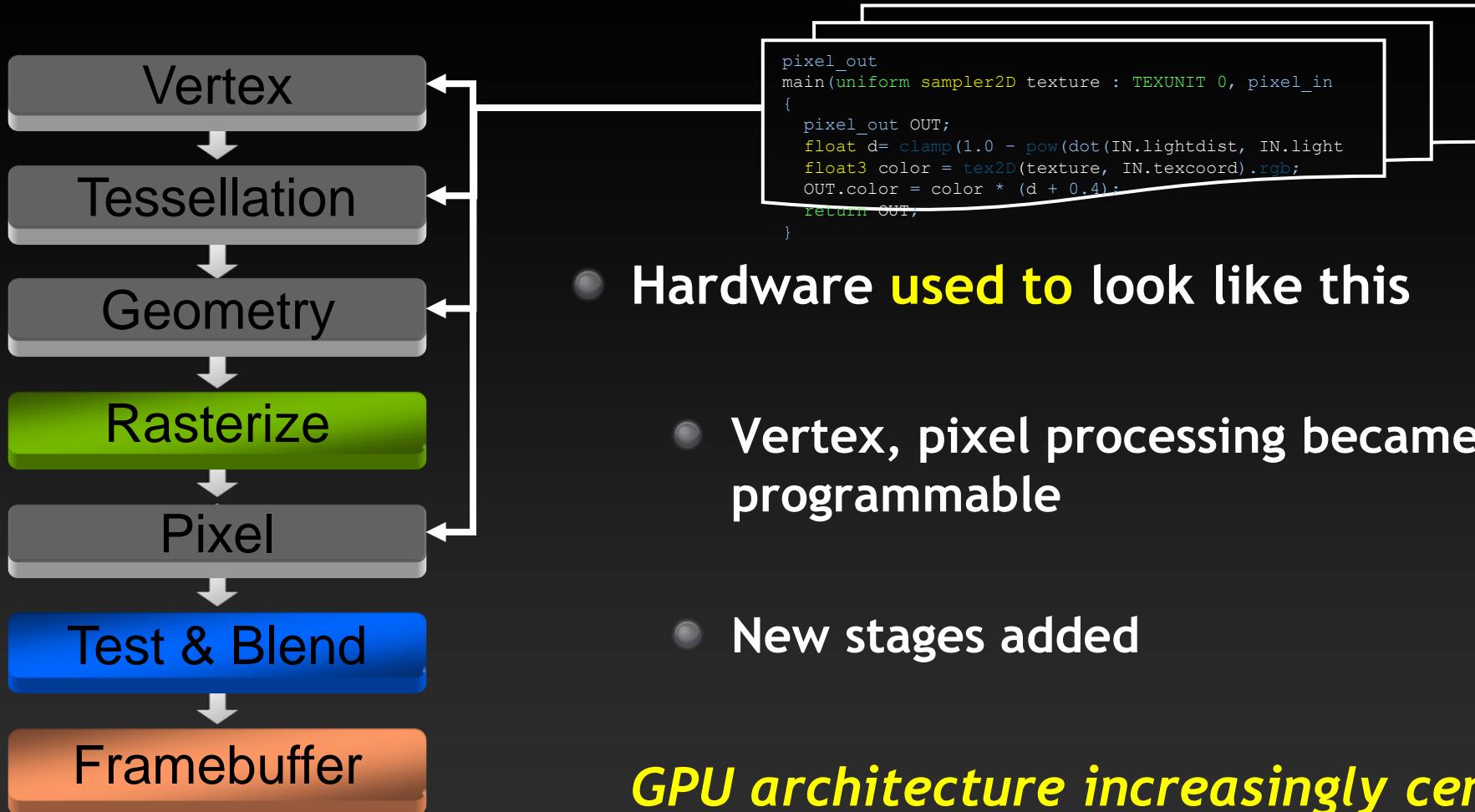


The Graphics Pipeline



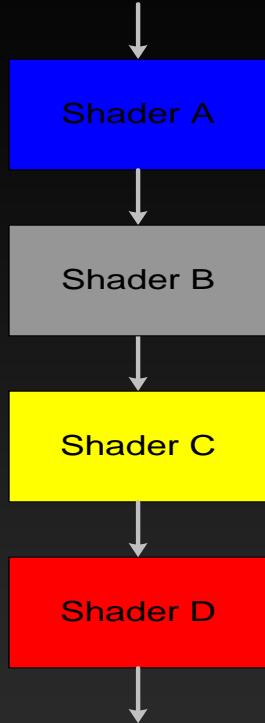
- **Hardware used to look like this**
- **Vertex, pixel processing became programmable**
- **New stages added**

The Graphics Pipeline

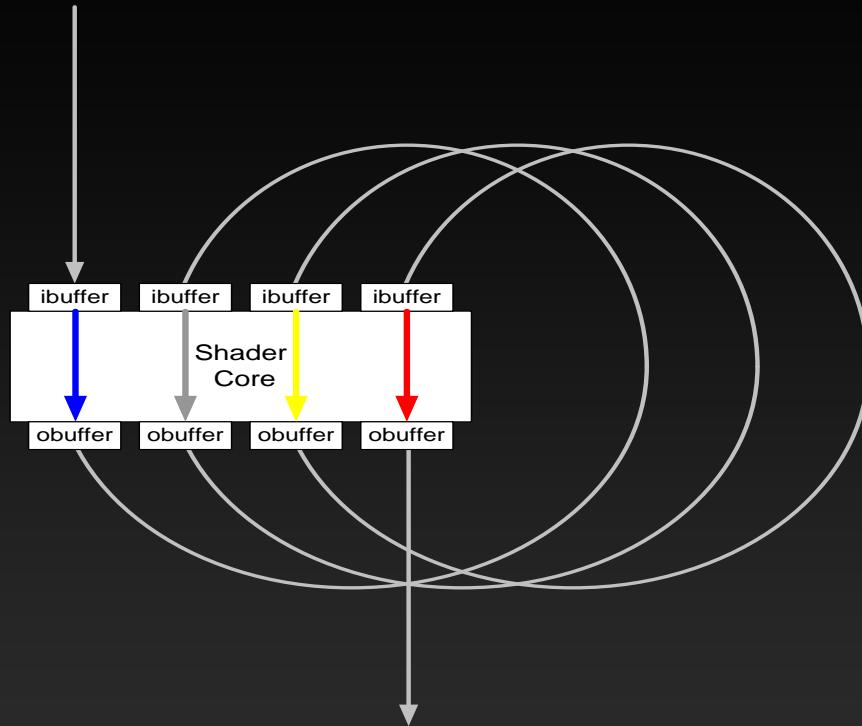


Modern GPUs: Unified Design

Discrete Design



Unified Design



Vertex shaders, pixel shaders, etc. become *threads*
running different programs on a flexible core

GeForce 8: First Fully Programmable GPU

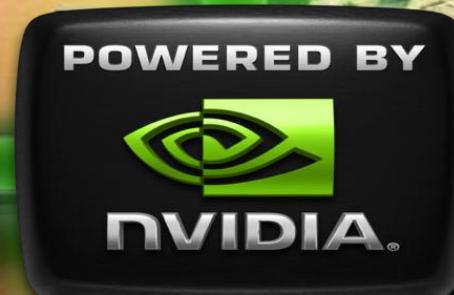


Stunning Graphics Realism

Lush, Rich Worlds



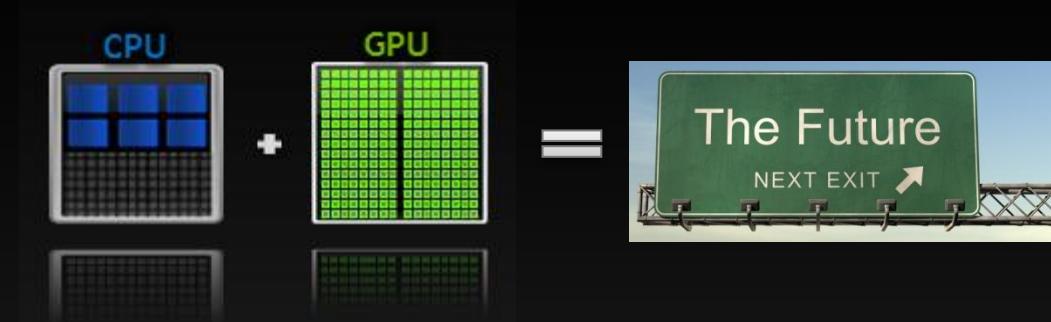
Crysis © 2006 Crytek / Electronic Arts



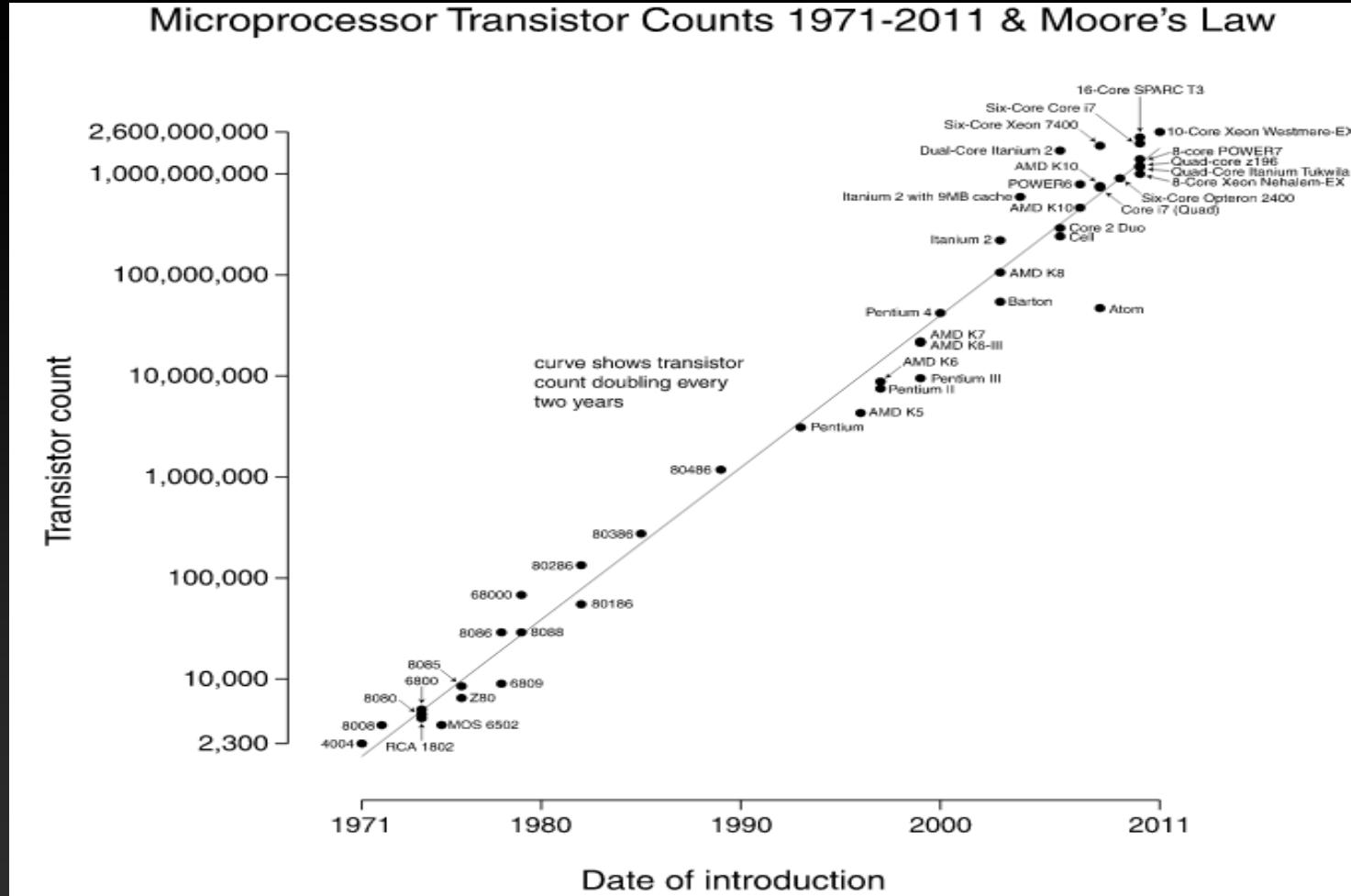
Incredible Physics Effects

Core of the Definitive Gaming Platform

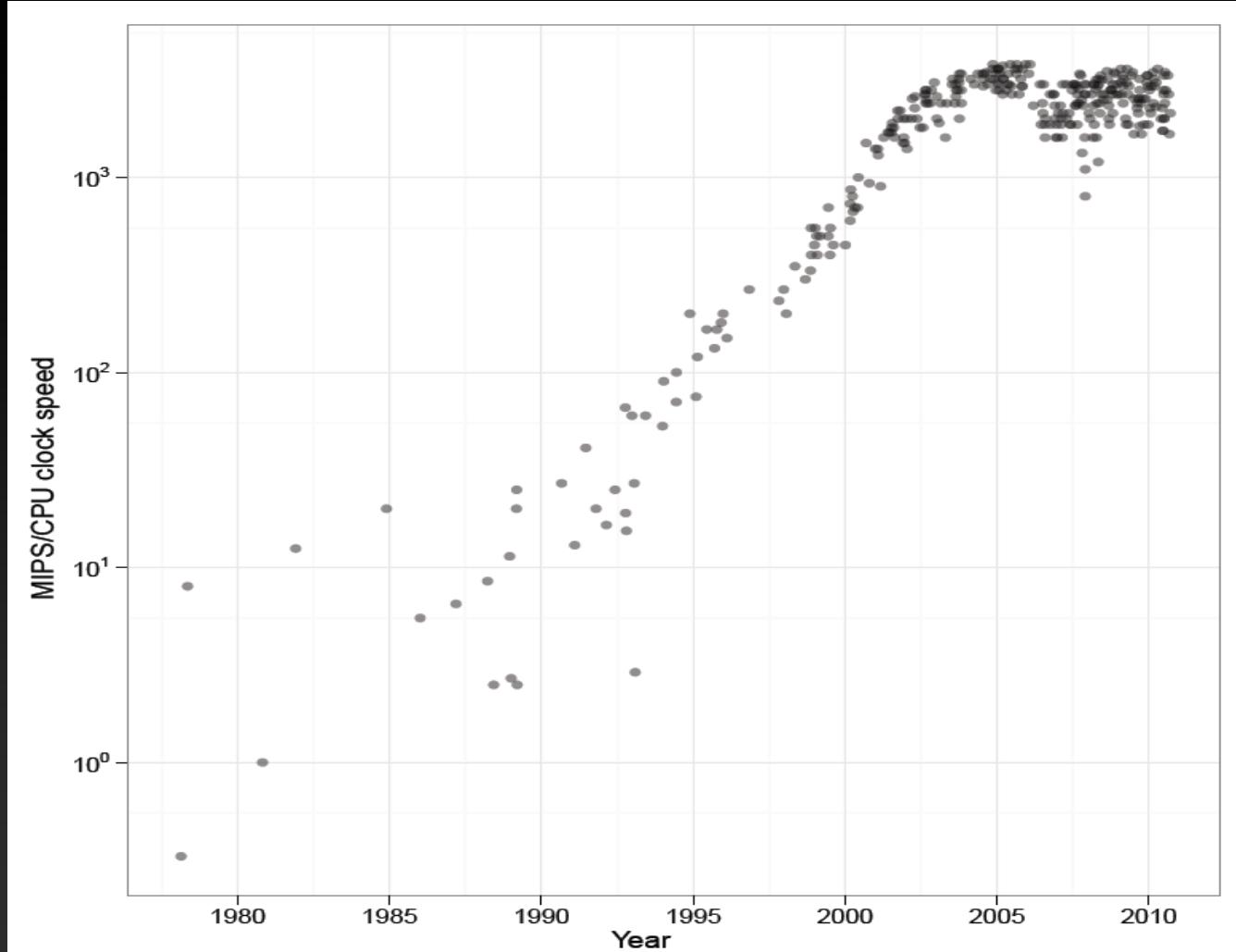
My Three Points



Moore's Law

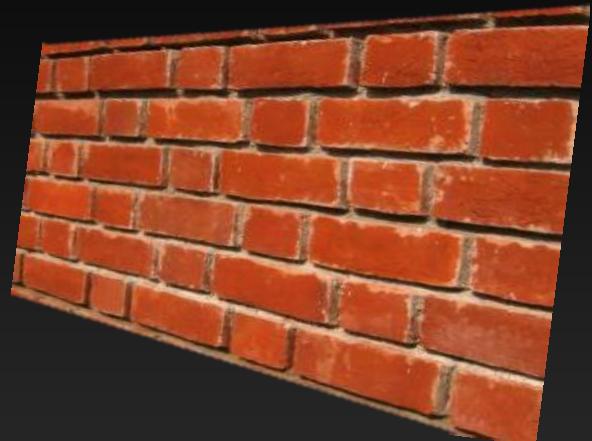


Moore's Law



What makes up the brick wall?

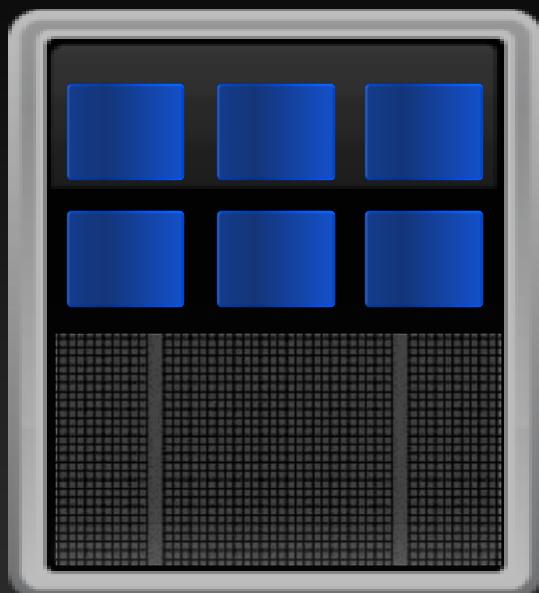
$$\begin{array}{c} \text{Power Wall} \\ + \\ \text{Memory Wall} \\ \hline \text{ILP Wall} \\ = \text{Brick Wall} \end{array}$$



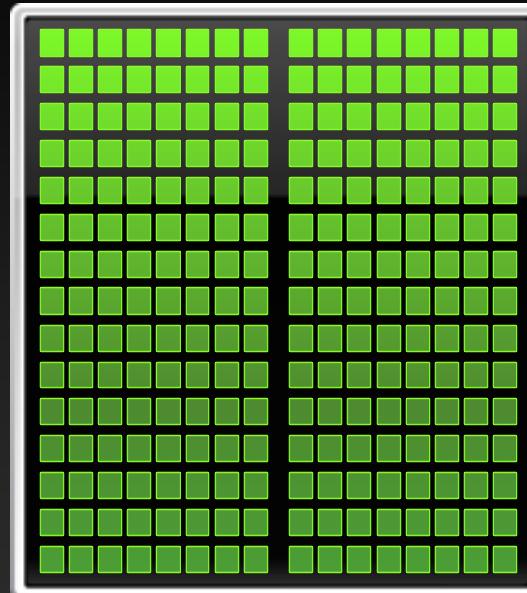
Heterogeneous Computing



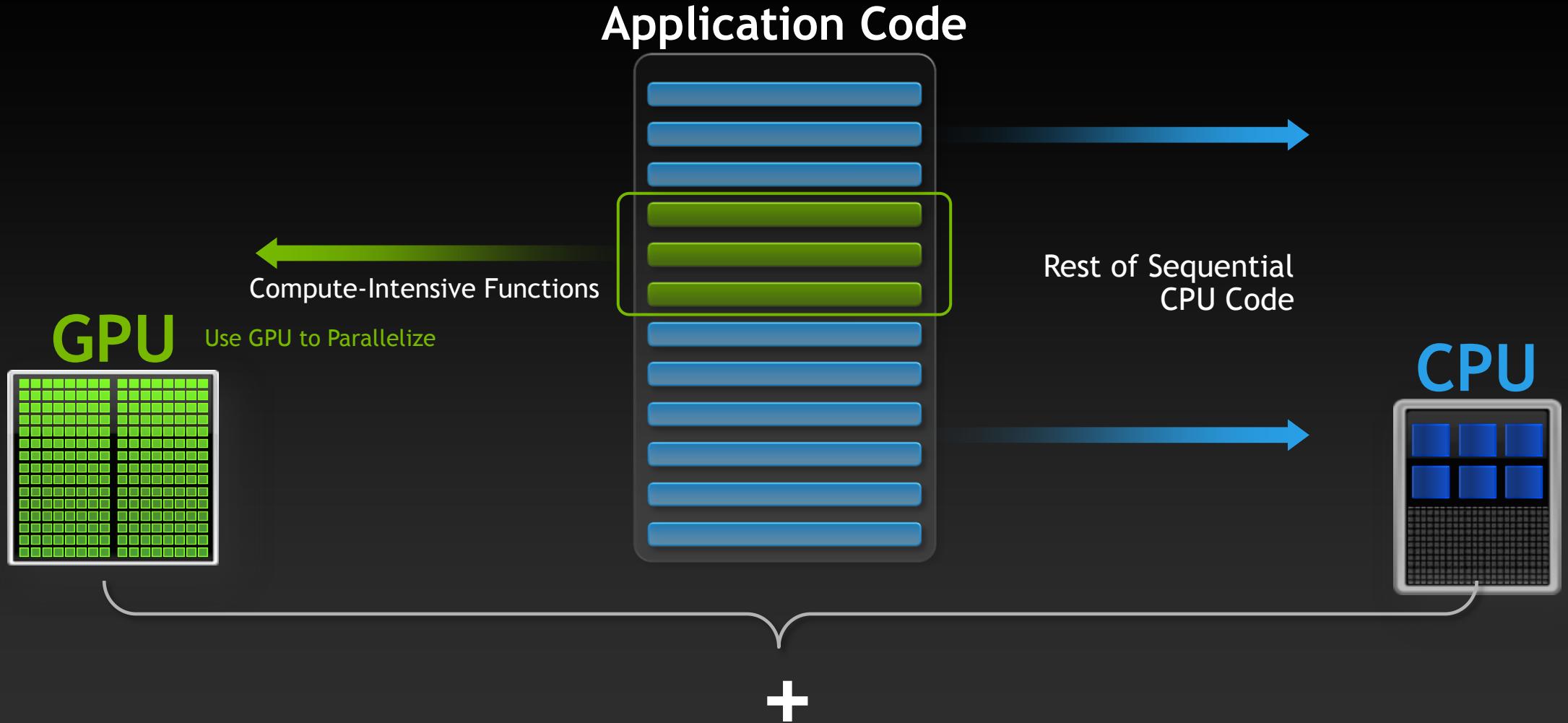
CPU



GPU

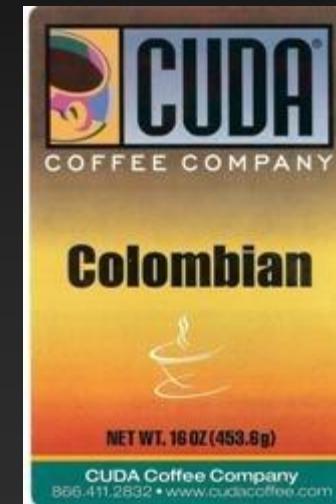


The basic idea



What is CUDA?

- Programming language?
- Compiler?
- Classic car?
- Beer?
- Wine?
- Coffee?



CUDA Parallel Computing Platform

www.nvidia.com/getcuda



Programming Approaches

Libraries

“Drop-in” Acceleration

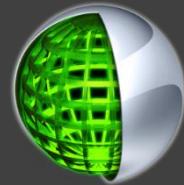
OpenACC Directives

Easily Accelerate Apps

Programming Languages

Maximum Flexibility

Development Environment



Nsight IDE

Linux, Mac and Windows
GPU Debugging and Profiling

CUDA-GDB debugger
NVIDIA Visual Profiler

Open Compiler Tool Chain



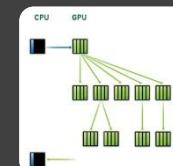
Enables compiling new languages to CUDA platform, and CUDA languages to other architectures

Hardware Capabilities

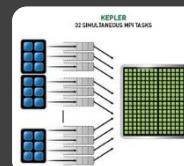
SMX



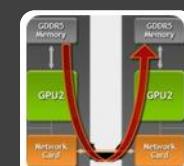
Dynamic Parallelism



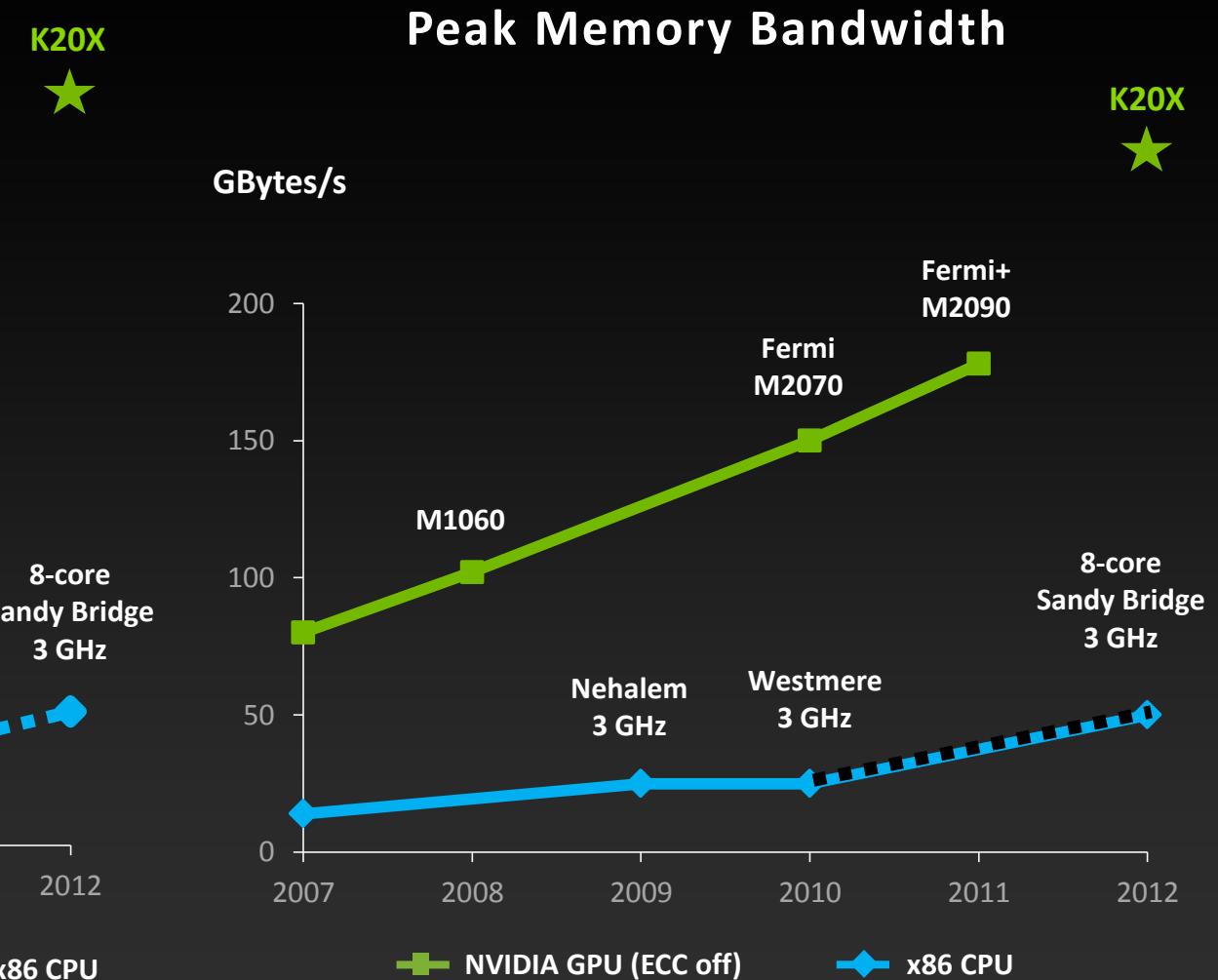
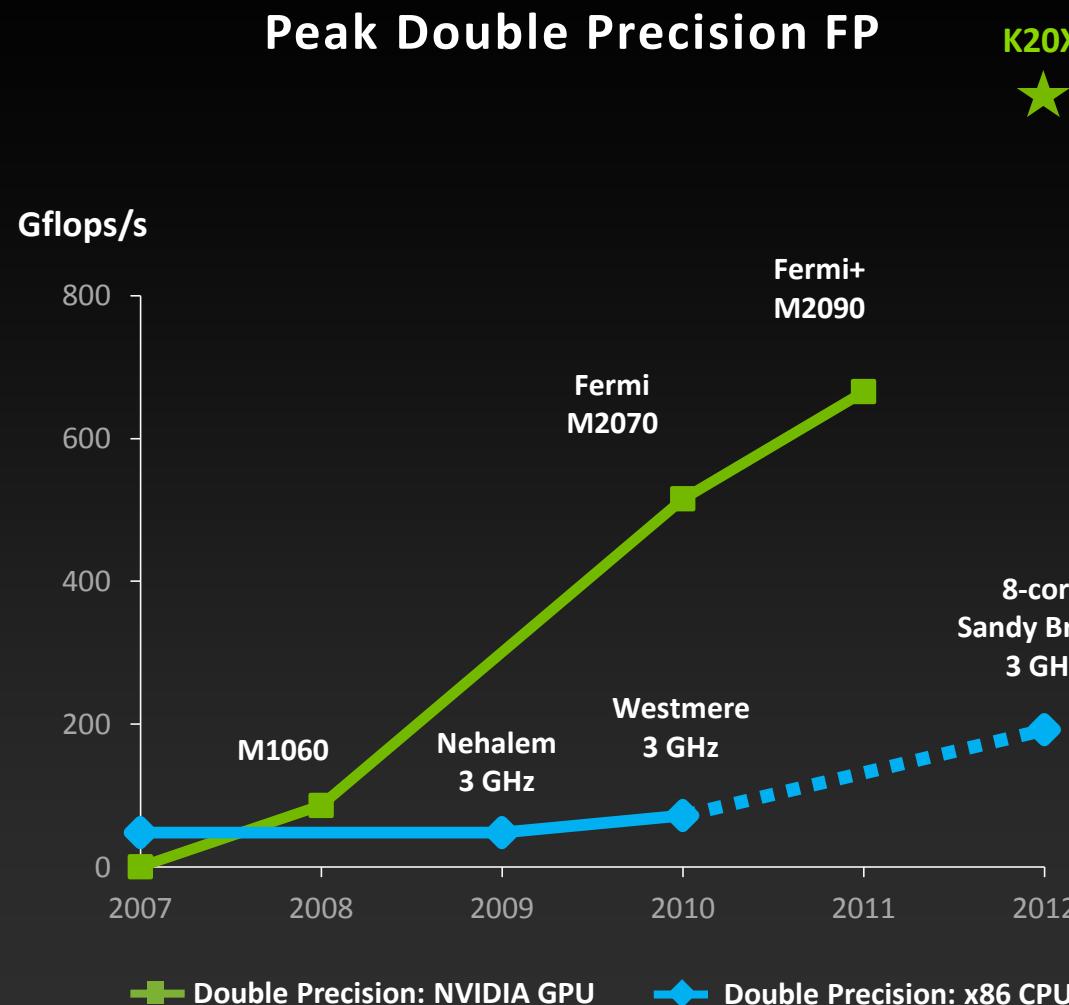
HyperQ



GPUDirect

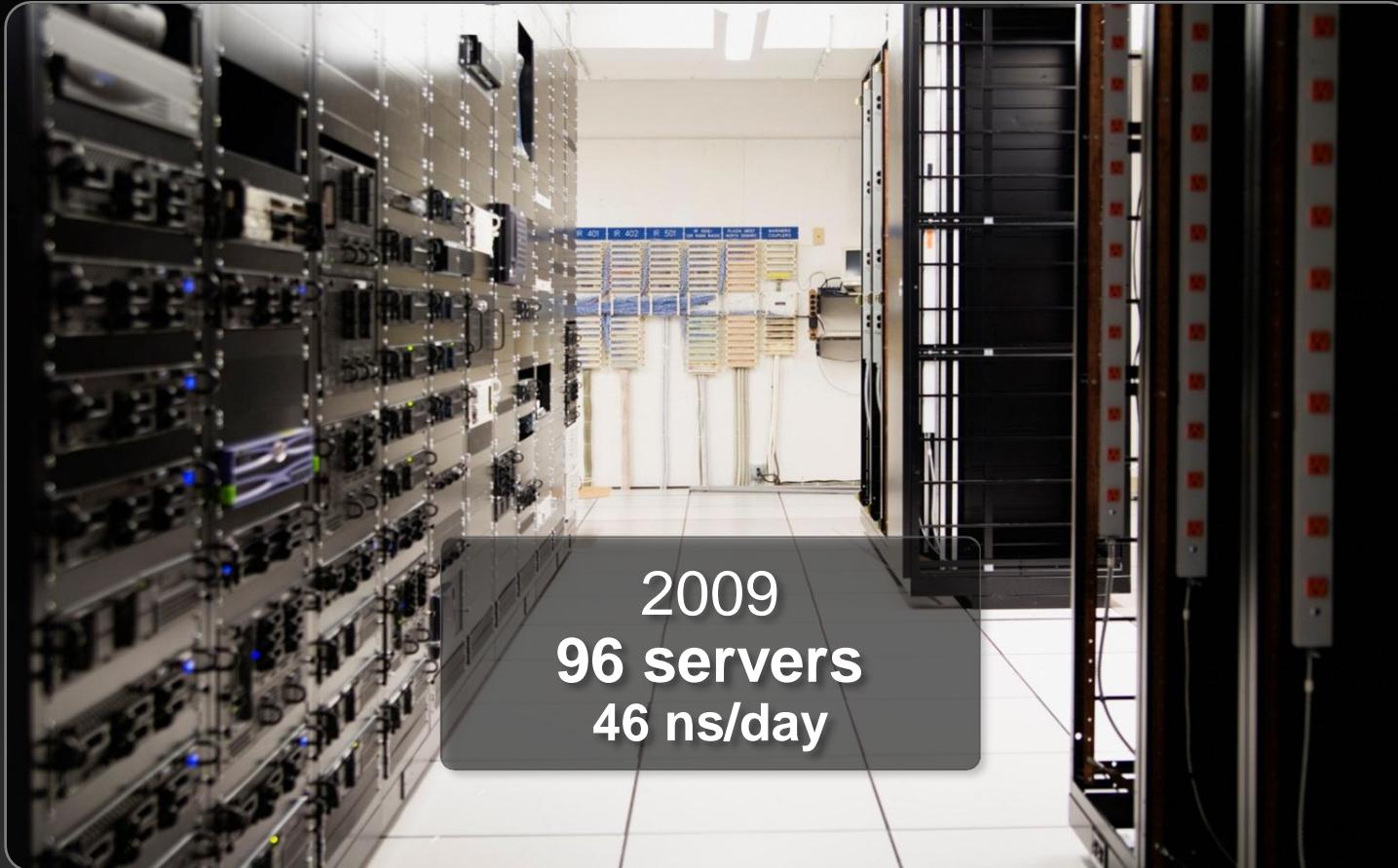


Why CUDA?

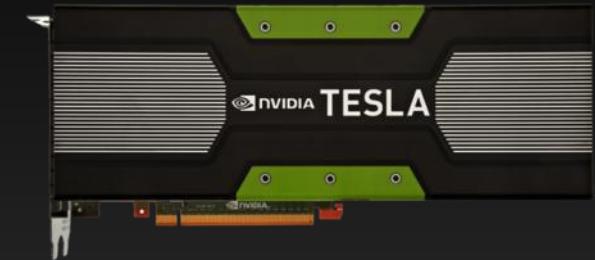


Kepler: Unprecedented Value to Science

AMBER-JAC NVE Benchmark



2009
96 servers
46 ns/day

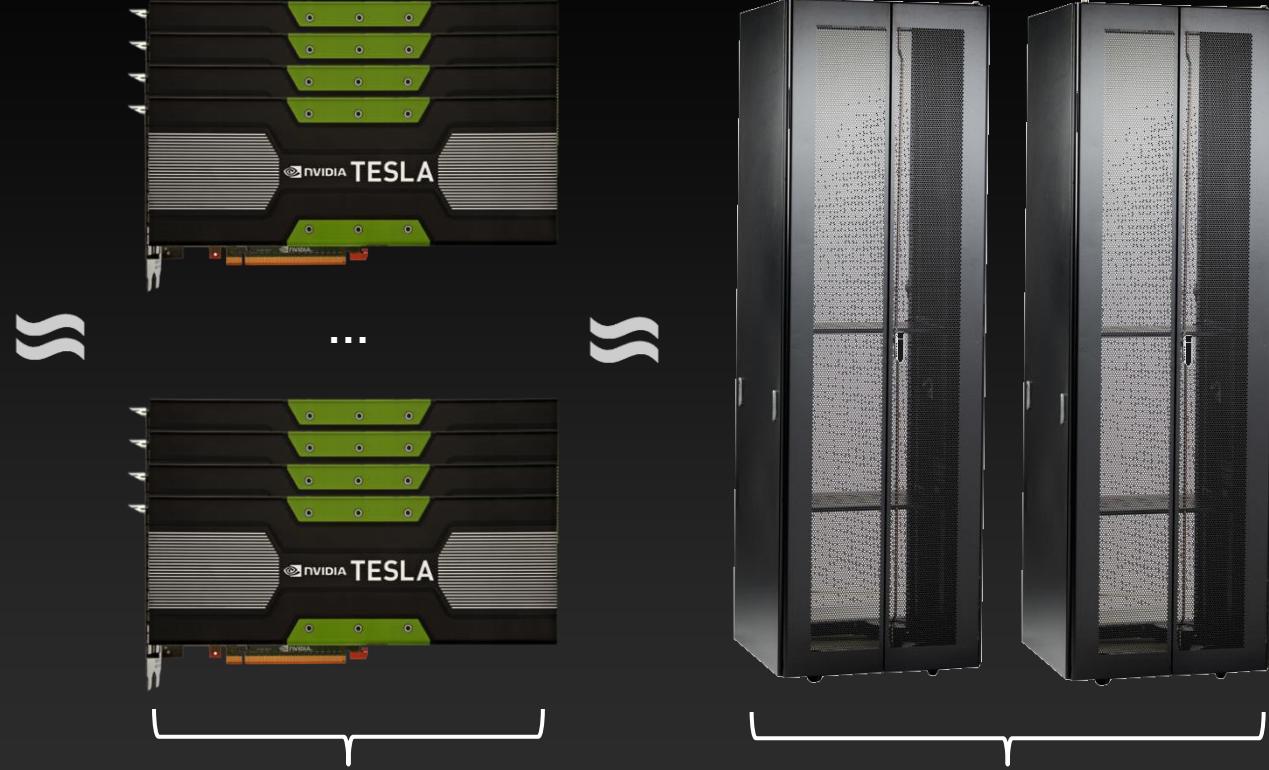


2012
1 Tesla K20 GPU
81 ns/day

JANUS vs. GPUs



JANUS: 184 Peak TFLOPs
(2010)

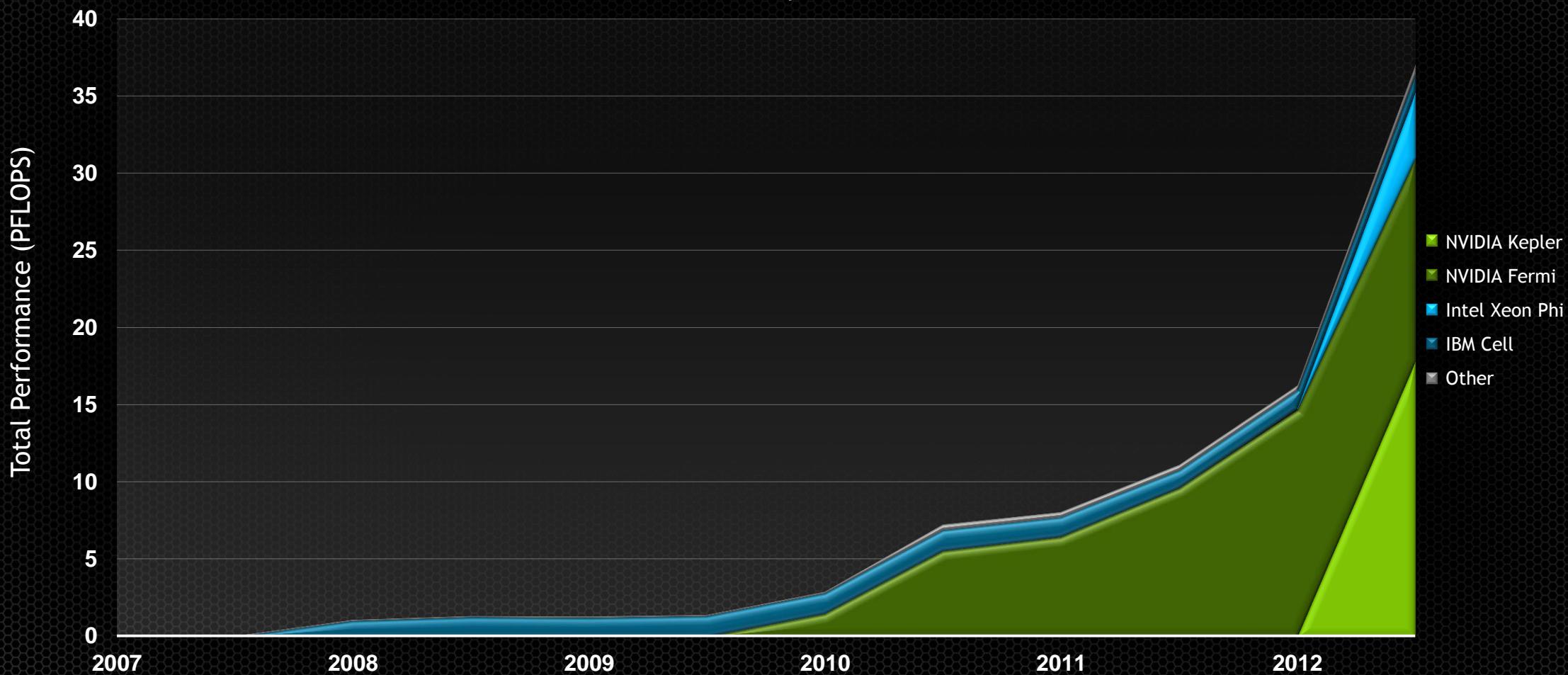


141 K20x GPUs
(2012)

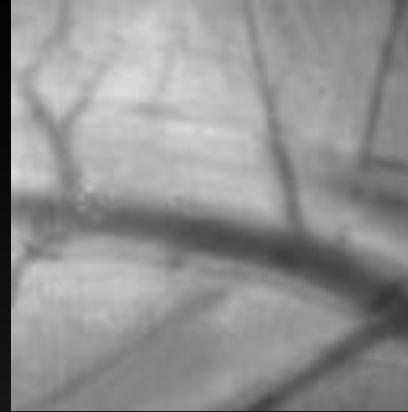
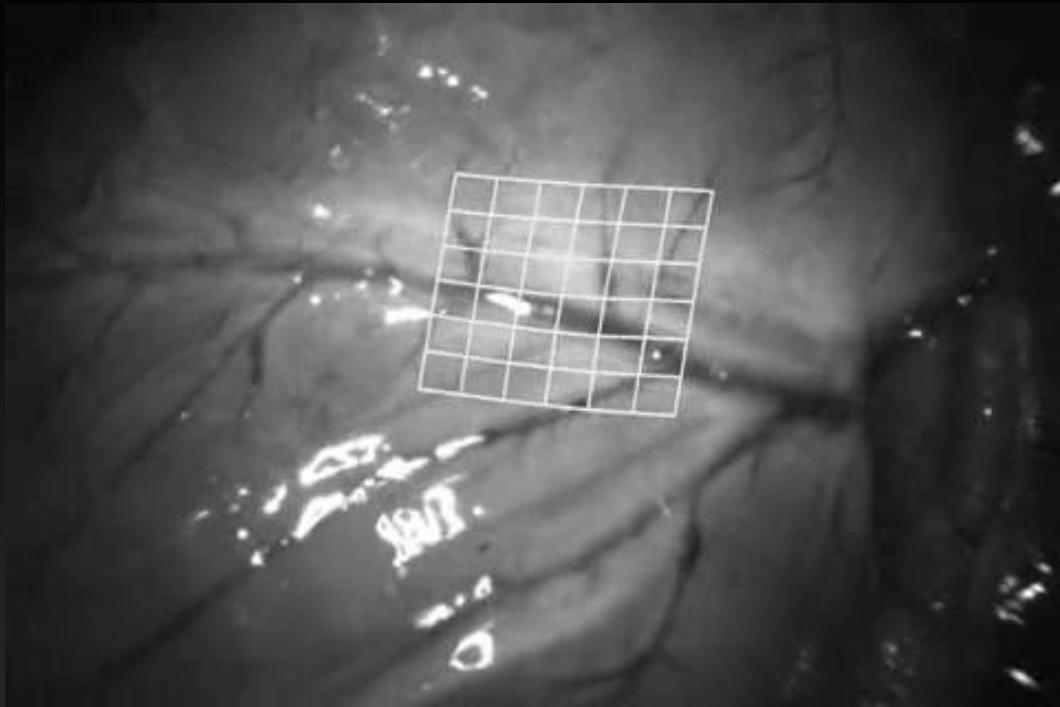
2 Server Racks

CUDA Accelerating

19% of FLOPS from GPU Systems



Operating on a Beating Heart



Only 2% of surgeons will operate on a beating heart

Patient stands to lose 1 point of IQ every 10 min with heart stopped

GPU enables real-time motion compensation to virtually stop beating heart for surgeons

Getting Started with CUDA



3 Ways to Accelerate Applications



Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

3 Ways to Accelerate Applications



Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

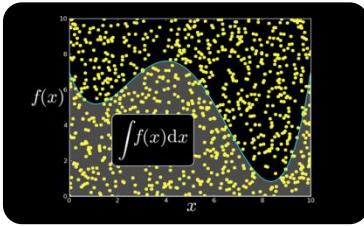


GPU Accelerated Libraries

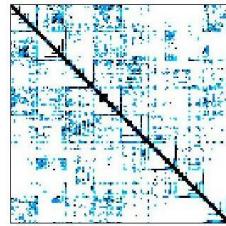
“Drop-in” Acceleration for your Applications



NVIDIA cuBLAS



NVIDIA cuRAND



NVIDIA cuSPARSE



NVIDIA NPP



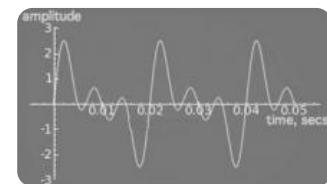
Vector Signal
Image Processing



GPU Accelerated
Linear Algebra



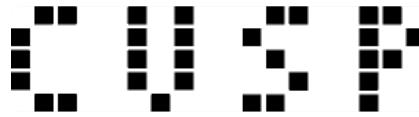
Matrix Algebra on
GPU and Multicore



NVIDIA cuFFT



ArrayFire Matrix
Computations



Sparse Linear
Algebra



C++ STL Features
for CUDA



3 Ways to Accelerate Applications



Applications

Libraries

“Drop-in”
Acceleration

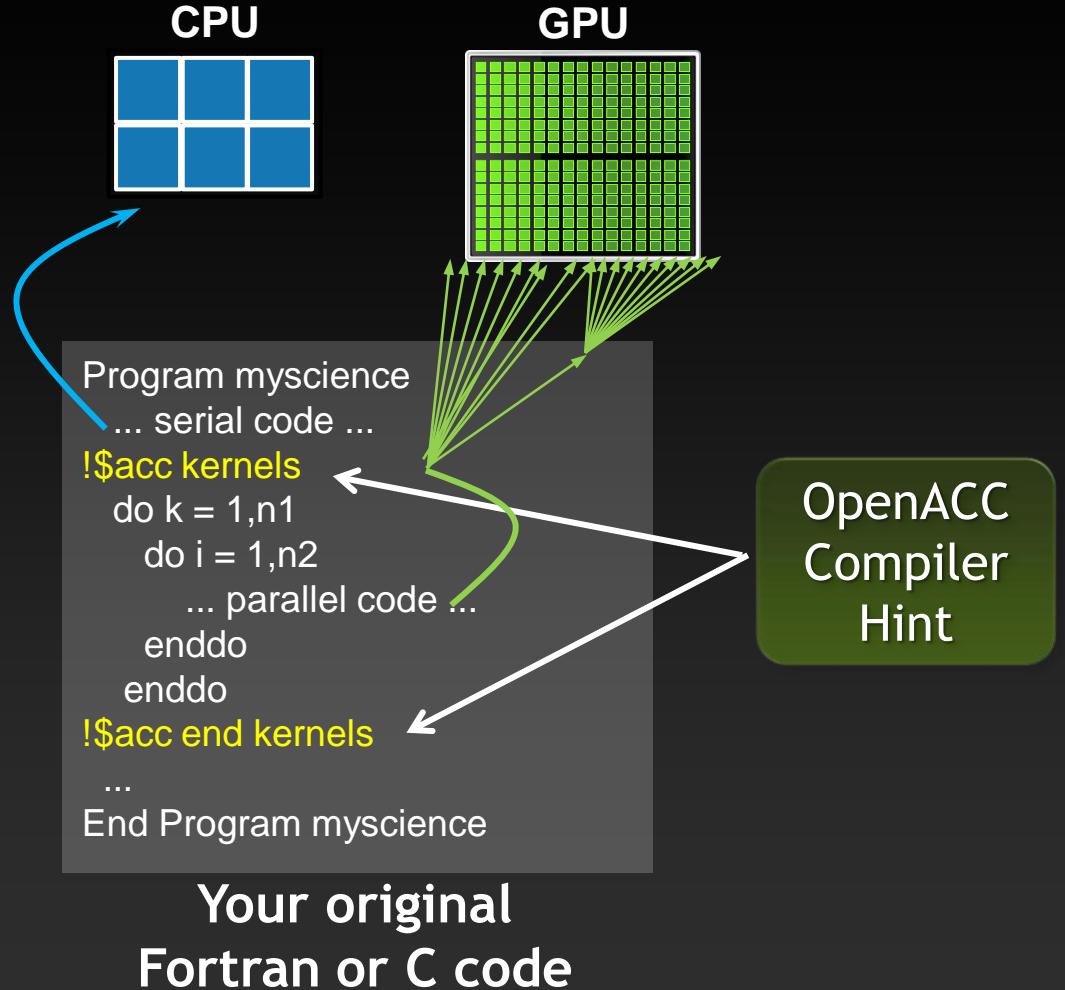
OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

OpenACC Directives



Simple Compiler hints

Compiler Parallelizes code

Works on many-core GPUs & multicore CPUs

3 Ways to Accelerate Applications



Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

GPU Programming Languages



Numerical analytics ►

MATLAB, Mathematica, LabVIEW

Fortran ►

OpenACC, CUDA Fortran

C ►

OpenACC, CUDA C

C++ ►

Thrust, CUDA C++

Python ►

PyCUDA, Copperhead, CUDA Python

F# ►

Alea.cubase

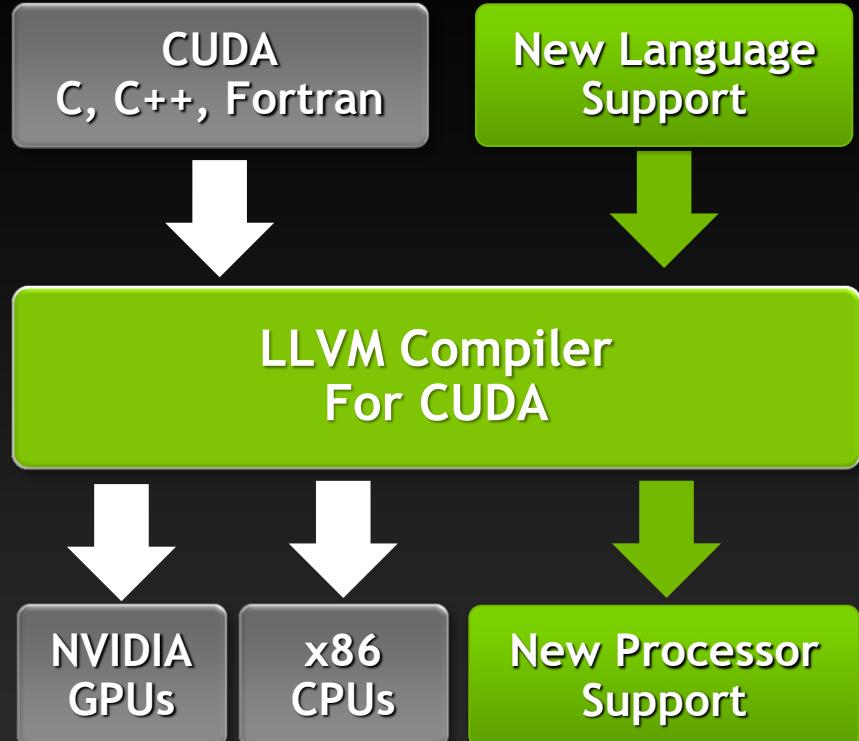
Enabling More Programming Languages



Developers want to build
front-ends for

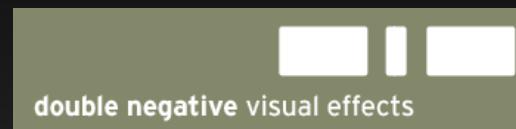
Python, Java, R, DSLs ...

Target other processors like
ARM, FPGAs, GPUs, x86 ...





Enabling More Programming Languages



Mozilla Rust

Halide (<http://halide-lang.org/>)

CUDA
C, C++, Fortran

New Language Support

LLVM Compiler
For CUDA

NVIDIA
GPUs

x86
CPUs

New Processor
Support





Rapid Development

Powerful Libraries



NumPy



SciPy.org



matplotlib

Large Community



Commercial Support



Is Python Fast Enough for HPC?



Python apps often implement
performance critical functions in C/C++.

PyCUDA



- Provides access to CUDA API from Python
 - Convenience classes for device data and code management
- Kernels must still be written in CUDA C
 - Can be dynamically generated and compiled on the fly by PyCUDA
- Open Source project



Copperhead

Standard Python

```
import numpy as np

def saxpy(a, x, y):
    return [a * xi + yi
            for xi, yi in zip(x, y)]

x = np.arange(2**20, dtype=np.float32)
y = np.arange(2**20, dtype=np.float32)

cpu_result = saxpy(2.0, x, y)
```

Copperhead: Parallel Python

```
from copperhead import *
import numpy as np

@cu
def saxpy(a, x, y):
    return [a * xi + yi
            for xi, yi in zip(x, y)]

x = np.arange(2**20, dtype=np.float32)
y = np.arange(2**20, dtype=np.float32)

with places.gpu0:
    gpu_result = saxpy(2.0, x, y)

with places.openmp:
    cpu_result = saxpy(2.0, x, y)
```

CU

Compile Python for Parallel Architectures



- Anaconda Accelerate from Continuum Analytics
 - NumbaPro array-oriented compiler for Python & NumPy
 - Compile for CPUs or GPUs (uses LLVM + NVIDIA Compiler SDK)
- Fast Development + Fast Execution: Ideal Combination



Free Academic
License

<http://continuum.io>

NumbaPro Demos



- **Mandelbrot iPython Notebook**
 - https://github.com/harrism/numbapro_examples (download)
 - <http://nbviewer.ipython.org/f5707335f40af9463c43> (view)
- **Monte Carlo iPython Notebook**
 - https://github.com/ContinuumIO/numbapro-examples/tree/master/monte_carlo_pricer (download)
 - <http://nbviewer.ipython.org/835a8ca39ced77fe751d> (view)

JANUS GPUs



GPU nodes

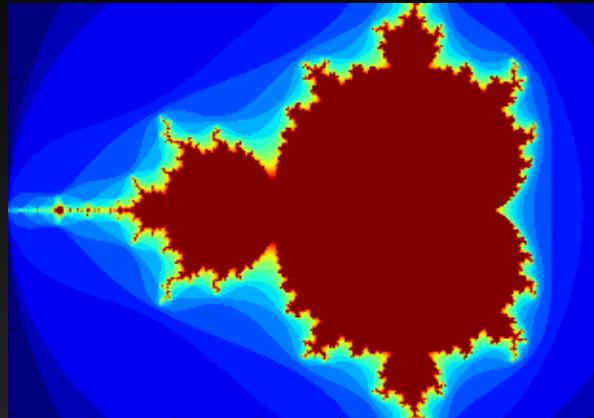
We currently provide two GPU-enabled nodes for visualization and GPU-accelerated applications. These nodes are accessed via the `crc-gpu` queue.

Their specifications are:

- 128 GB RAM
- 2 1-TB local disks
- 12 physical cores (24 with hyperthreading), Intel Xeon X5660 processors (2.8 GHz)
- 1 NVIDIA Tesla M2070 GPU

The `crc-gpu` queue has a maximum walltime of 4 hours.

CUDA Python



```
@cuda.jit(restype=uint32, argtypes=[f8, f8, uint32], device=True)
def mandel(x, y, max_iters):
    c = complex(x, y)
    z = 0.0j
    for i in range(max_iters):
        z = z*z + c
        if (z.real*z.real + z.imag*z.imag) >= 4:
            return i
    return max_iters

@cuda.jit(argtypes=[uint8[:, :], f8, f8, f8, f8, uint32])
def mandel_kernel(img, xmin, xmax, ymin, ymax, iters):
    x, y = cuda.grid(2)
    if x < img.shape[0] and y < img.shape[1]:
        img[y, x] = mandel(min_x+x*((max_x-min_x)/img.shape[0]),
                            min_y+y*((max_y-min_y)/img.shape[1]), iters)

gimage = np.zeros((1024, 1024), dtype = np.uint8)
d_image = cuda.to_device(gimage)
mandel_kernel[(32,32), (32,32)](d_image, -2.0, 1.0, -1.0, 1.0, 20)
d_image.to_host()
```

CUDA Programming, Python Syntax

	1024 ² Mandelbrot Time	Speedup v. Pure Python
Pure Python	4.85s	--
NumbaPro (CPU)	0.11s	44x
CUDA Python (K20)	.004s	1221x

CUDA Breeds Academic Success



SET EDITION: U.S. | INTERNATIONAL | MEXICO | ARABIC
TV: CNN | CNN en Español | HLN

Sign up | Log in

CNN Tech

SEARCH

POWERED BY Google

Home Video NewsPulse U.S. World Politics Justice Entertainment Tech Health Living Travel Opinion iReport Money Sports

PopTech: 5 fascinating people you've never heard of

By John D. Sutter, CNN

updated 8:32 AM EST, Mon October 24, 2011 | Filed under: [Innovations](#)

[Twitter](#) [Share](#) [Email](#) [Save](#) [Print](#)

[Recommend](#)

326 people recommend this. Be the first of your friends.

Iain Couzин uses Xbox Kinect cameras and computer vision to study swarm behavior in creatures like locusts, fish and humans.

STORY HIGHLIGHTS

• PopTech is a yearly conference that focuses on tech and social change

Camden, Maine (CNN) -- There are no rock stars at PopTech, no household names. But this annual conference in coastal Maine is a hub for super-smart people, a chance to get a look into ideas and technologies that soon will change the world.

PAID FOR BY OBAMA FOR AMERICA

ADVERTISING

[TECH: NEWSPULSE]

Most popular Tech stories right now

A new, high-definition iPad from Apple [progress bar]

Review: New iPad neither dud nor revolution [progress bar]

The "new iPad" doesn't have a new name [progress bar]



The image shows the front page of The New York Times Science Times section. At the top, there's a large, artistic illustration of many fish swimming in a dark blue, textured ocean. Below the illustration, the word "Science Times" is written in a bold, serif font. To the right of the title, there's a small caption that reads "Illustration: ROBERT DUNCAN". On the far left, there's a small box containing the text "100% Recycled Paper". On the far right, there's a small box containing the text "National Edition". The main headline at the bottom of the page is "From Ants to People, an Instinct to Swarm". There are several columns of text and some smaller illustrations or maps to the right of the main headline.

Facial Recognition with CUDA



- NICTA: National Information and Communication Technology Australia
 - Specialize in security research

With fricken laserbeams!



- Created by Intellectual Ventures to help fight malaria in third world countries
- Image detection and targeting is done with NVIDIA GPUs



Links to get started

- Get CUDA: www.nvidia.com/getcuda
- Anaconda Accelerate: store.continuum.io/cshop/accelerate/
- Programming Guide/Best Practices...
 - docs.nvidia.com
- Questions:
 - NVIDIA Developer forums devtalk.nvidia.com
 - Search or ask on www.stackoverflow.com/tags/cuda
- General: www.nvidia.com/cudazone

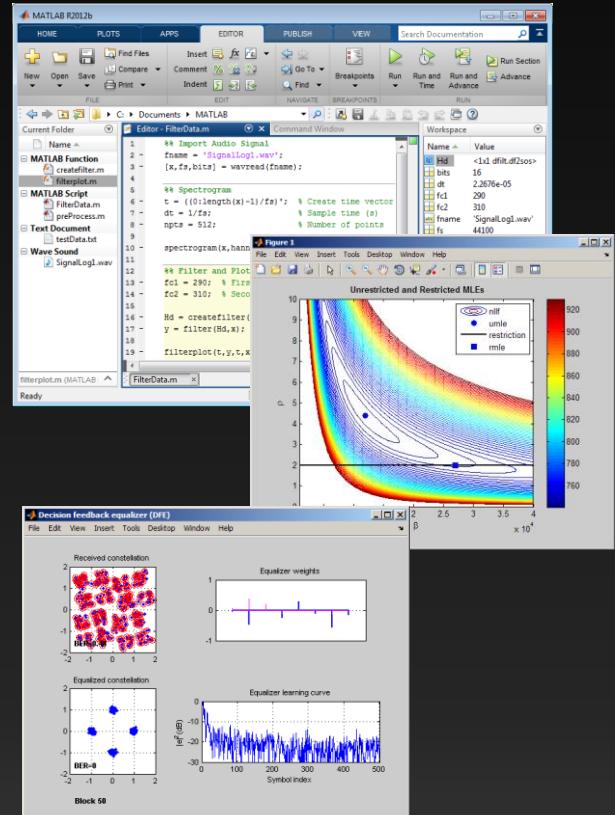
Parallel Computing with MATLAB

On-ramp to GPU Computing

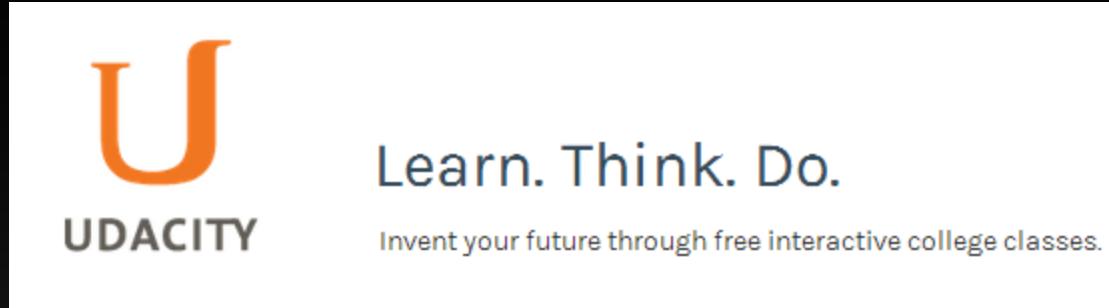


MATLAB

- Over 200 of the most popular MATLAB functions on GPUs
Including:
 - Random number generation
 - FFT
 - Matrix multiplications
 - Solvers
 - Convolutions
 - Min/max
 - SVD
 - Cholesky and LU factorization
- MATLAB Compiler support (GPU acceleration in stand-alone applications)
- GPU features in
 - Communication Systems Toolbox
 - Phased Array System Toolbox
 - Signal Processing Toolbox
 - Neural Network Toolbox



Udacity Parallel Programming Course



The image shows a snippet of the Udacity website. It features the orange 'U' logo and the word 'UDACITY' in a sans-serif font. Below the logo is the tagline 'Learn. Think. Do.' followed by the subtitle 'Invent your future through free interactive college classes.'



20,000+ Students Already Taking
the course!

IT'S FREE!

[Udacity.com](https://www.udacity.com)

