

# Languages, APIs and Development Tools for GPU Computing

Phillip Miller  
Director, Software Product Management  
Professional Solutions Group



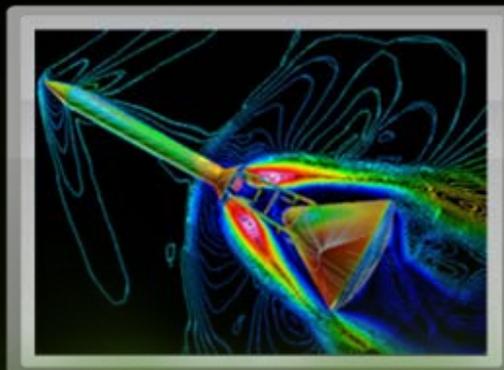
# CUDA Architecture

## Available on All Modern NVIDIA GPUs

**GeForce®**  
Entertainment



**Tesla™**  
High Performance Computing



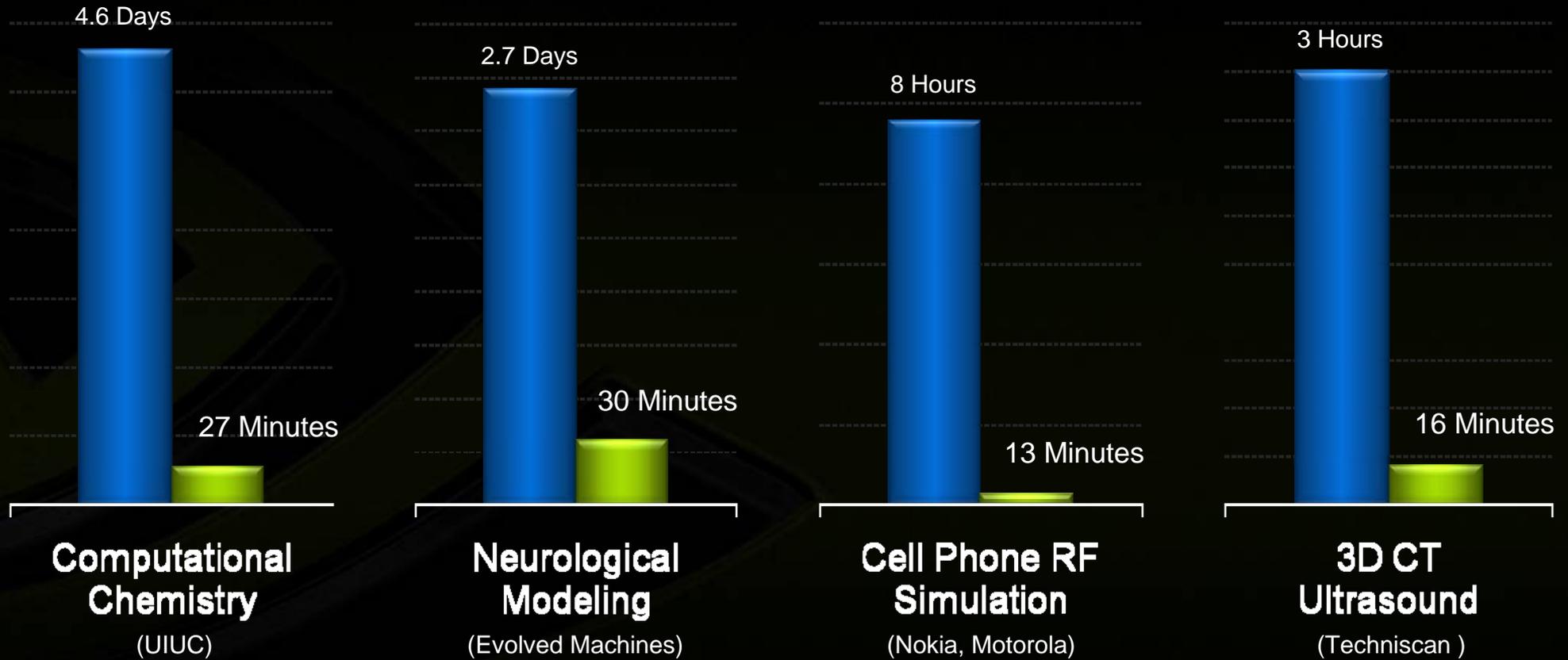
**Quadro®**  
Design & Creation



**GPU**

**Over 150 Million CUDA-enabled GPUs Installed**

# GPUs are Accelerating Time to Discovery



# Huge Speed-Ups Across Many Fields

Algorithm	Field	Speedup
2-Electron Repulsion Integral	Quantum Chemistry	<b>130X</b>
Lattice Boltzmann	CFD	<b>123X</b>
Euler Solver	CFD	<b>16X</b>
GROMACS	Molecular Dynamics	<b>137X</b>
Lattice QCD	Physics	<b>30X</b>
Multifrontal Solver	FEA	<b>20X</b>
nbody	Astrophysics	<b>100X</b>
Simultaneous Iterative Reconstruction Technique	Computed Tomography	<b>32X</b>

# GPU Computing Overview

## Broad Adoption

- Over 150,000,000 installed CUDA-Architecture GPUs
- Over 90,000 GPU Computing Developers (9/09)
- Windows, Linux and MacOS Platforms supported
- GPU Computing spans HPC to Consumer
- 250+ Universities teaching GPU Computing on the CUDA Architecture

## GPU Computing Applications

### CUDA C/C++

- Over 90,000 developers
- Running in Production since 2008
- SDK + Libs + Visual Profiler and Debugger

### OpenCL

- 1<sup>st</sup> GPU demo
- Shipped 1<sup>st</sup> OpenCL Conformant Driver
- Public Availability (Since April)

### Direct Compute

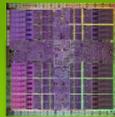
- Microsoft API for GPU Computing
- Supports all CUDA-Architecture GPUs (DX10 and DX11)

### Fortran

- PGI Accelerator
- PGI CUDA Fortran
- NOAA Fortran bindings
- FLAGON

### Python, Java, .NET, ...

- PyCUDA
- jCUDA
- CUDA.NET
- OpenCL.NET



## NVIDIA GPU

with the CUDA Parallel Computing Architecture

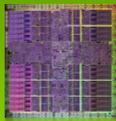
# GPU Computing Application Development

**Your GPU Computing Application**

**Application Acceleration Engines (AXEs)**  
Middleware, Modules & Plug-ins

**Foundation Libraries**  
Low-level Functional Libraries

**Development Environment**  
Languages, Device APIs, Compilers, Debuggers, Profilers, etc.

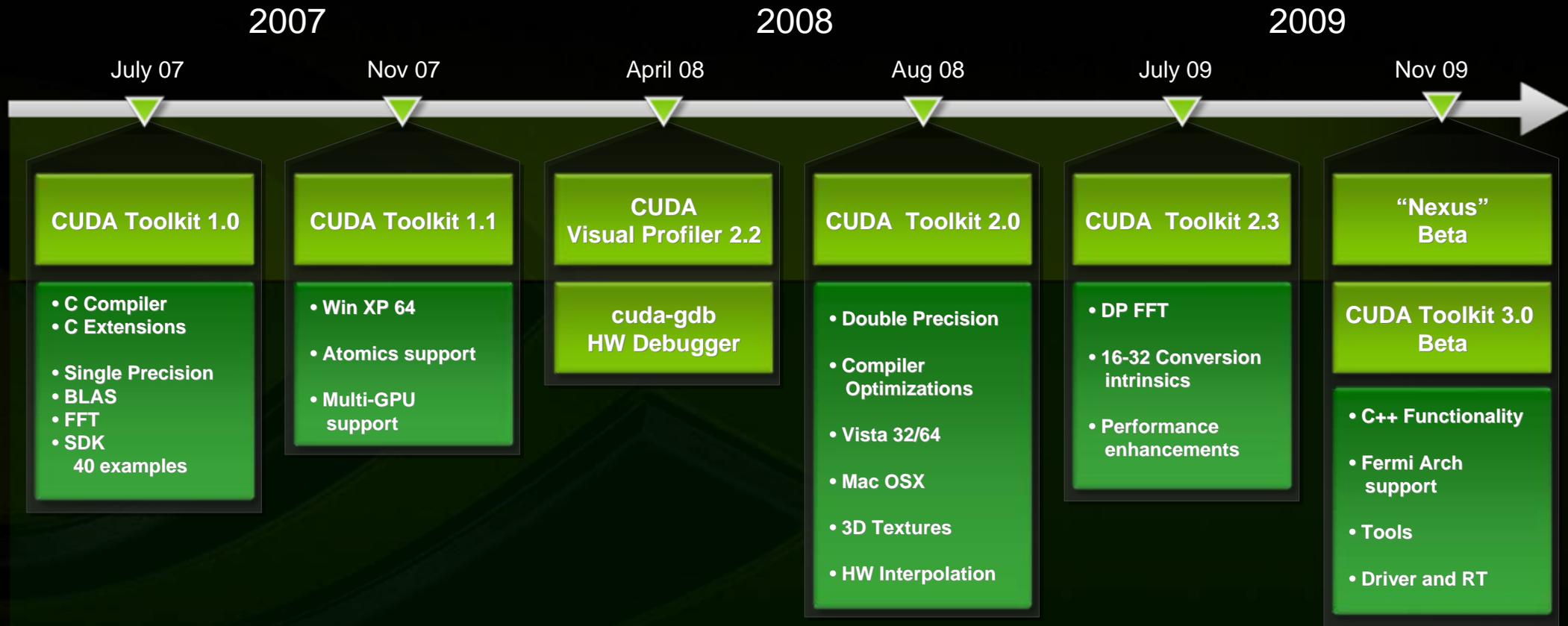


**CUDA Architecture**



# Languages & APIs

# CUDA C/C++ Update



# CUDA C/C++ Update

## CUDA C++ Language Functionality

- C++ Class Inheritance
- C++ Template Inheritance

(Productivity)

## Fermi Architecture Support

- Native 64-bit GPU support
- Generic Address space
- Dual DMA support
- ECC reporting
- Concurrent Kernel Execution

Architecture supported in SW now

## Tools

- “Nexus” Visual Studio IDE
- Debugger support for CUDA Driver API
- ELF support (cubin format deprecated)
- CUDA Memory Checker (cuda-gdb feature)
- Fermi: Fermi HW Debugger support in cuda-gdb
- Fermi: Fermi HW Profiler support in CUDA Visual Profiler

## CUDA Driver & CUDA C Runtime

- CUDA Driver / Runtime Buffer interoperability (Stateless CUDART)
- Separate emulation runtime
- Unified interoperability API for Direct3D and OpenGL
- OpenGL texture interoperability
- Fermi: Direct3D 11 interoperability

2009

July 09

Nov 09

“Nexus”  
Beta

CUDA Toolkit 3.0  
Beta

- C++ Functionality
- Fermi Arch support
- Tools
- Driver and RT

# Fortran Language Solutions

- **3 PGI Accelerators**
  - High-level *implicit* programming model, similar to OpenMP
  - Auto-parallelizing compiler
- **4 PGI CUDA Fortran Compiler**
  - High-level *explicit* programming model, similar to CUDA C Runtime
- **1 NOAA F2C-ACC**
  - Converts Fortran codes to CUDA C
  - Some hand-optimization expected
- **2 FLAGON**
  - Fortran 95 Library for GPU Numerics
  - Includes support for cuBLAS, cuFFT, CUDPP, etc.

# OpenCL

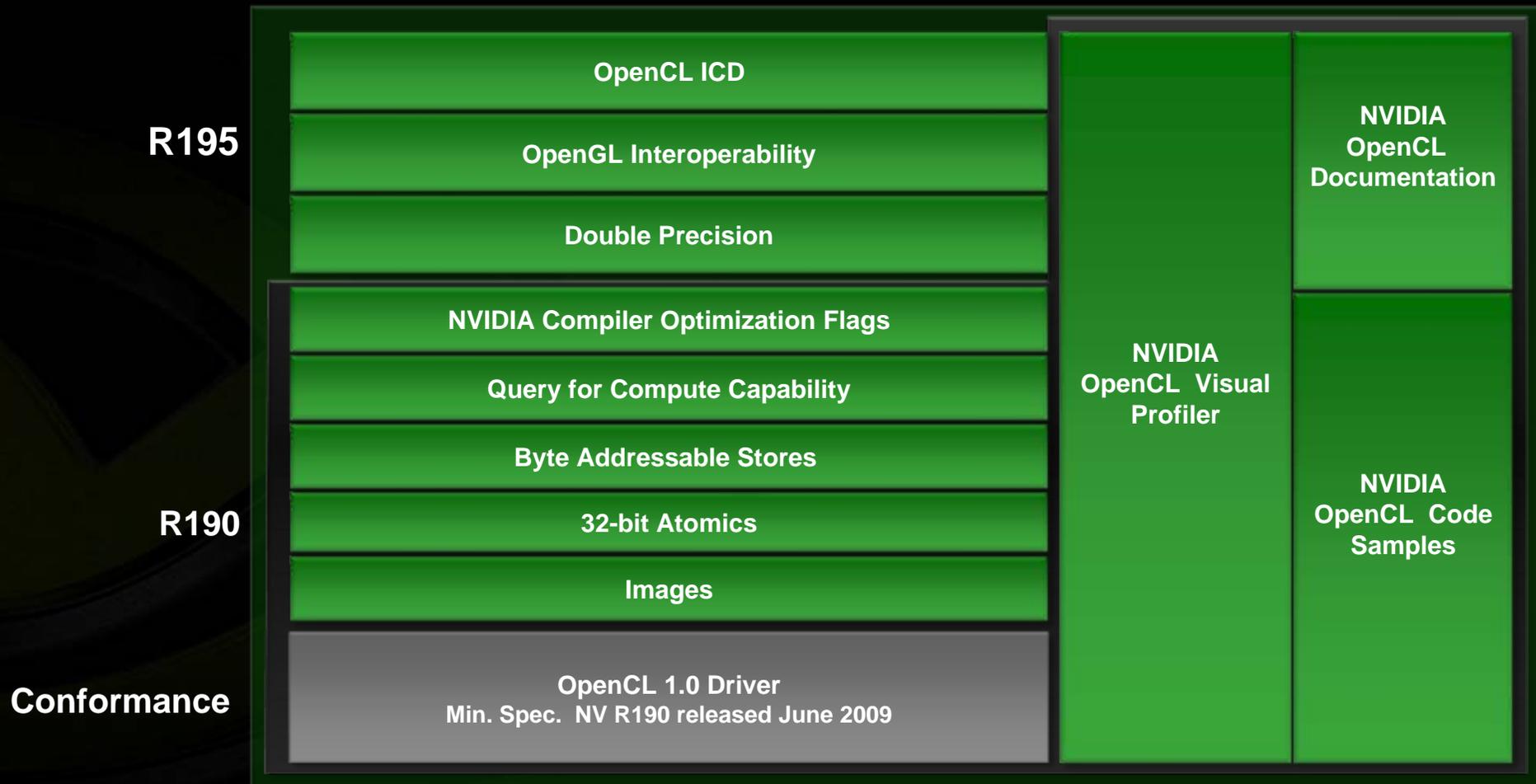
- **Cross-vendor open standard**
  - Managed by the Khronos Group
- **Low-level API for device management and launching kernels**
  - Close-to-the-metal programming interface
  - JIT compilation of kernel programs
- **C-based language for compute kernels**
  - Kernels must be optimized for each processor architecture



<http://www.khronos.org/opencvl>

NVIDIA released the first OpenCL v1.0 conformant driver for Windows and Linux to thousands of developers in June 2009

# NVIDIA OpenCL Support



# DirectCompute

- **Microsoft standard for all GPU vendors**
  - Released with DirectX® 11 / Windows 7
  - Runs on all 150M+ CUDA-enabled DirectX 10 class GPUs and later
- **Low-level API for device management and launching kernels**
  - Good integration with other DirectX APIs
- **Defines HLSL-based language for compute shaders**
  - Kernels must be optimized for each processor architecture

# Language & APIs for GPU Computing

Solution	Approach
CUDA C / C++	Language Integration Device-Level API
PGI Accelerator PGI CUDA Fortran	Auto Parallelizing Compiler Language Integration
CAPS HMPP	Auto Parallelizing Compiler
OpenCL	Device-Level API
DirectCompute	Device-Level API
PyCUDA	API Bindings
CUDA.NET, OpenCL.NET, jCUDA	API Bindings
...	...



# Development Tools

# NVIDIA Developer Resources

## DEVELOPMENT TOOLS

### CUDA Toolkit

Complete GPU computing development kit

### Visual Profiler

GPU hardware profiler for CUDA C and OpenGL

### Nexus

Development environment with Visual Studio integration [ *beta* ]

### NVPerfKit

OpenGL|D3D performance tools

### FX Composer

Shader Authoring IDE



## SDKs AND CODE SAMPLES

### GPU Computing SDK

CUDA C, OpenCL, DirectCompute programming guide and code samples

### Graphics SDK

DirectX & OpenGL code samples

### PhysX SDK

Complete game physics solution

### OpenAutomate

Test automation SDK



## VIDEO LIBRARIES

### Video Decode Acceleration

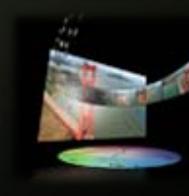
NVCUVID  
DXVA  
Win7 MFT

### Video Encode Acceleration

NVCUVENC  
Win7 MFT

### Post Processing

Noise reduction / De-interlace/  
Polyphase scaling / Color process



## ENGINES & LIBRARIES

### NPP Image Libraries

Performance primitives for imaging

### Numeric Libraries

cuFFT, cuLA, cuBLAS

### App Acceleration Engines

Optimized software modules for GPU acceleration

### Shader Library

Shader and post processing

### Optimization Guides

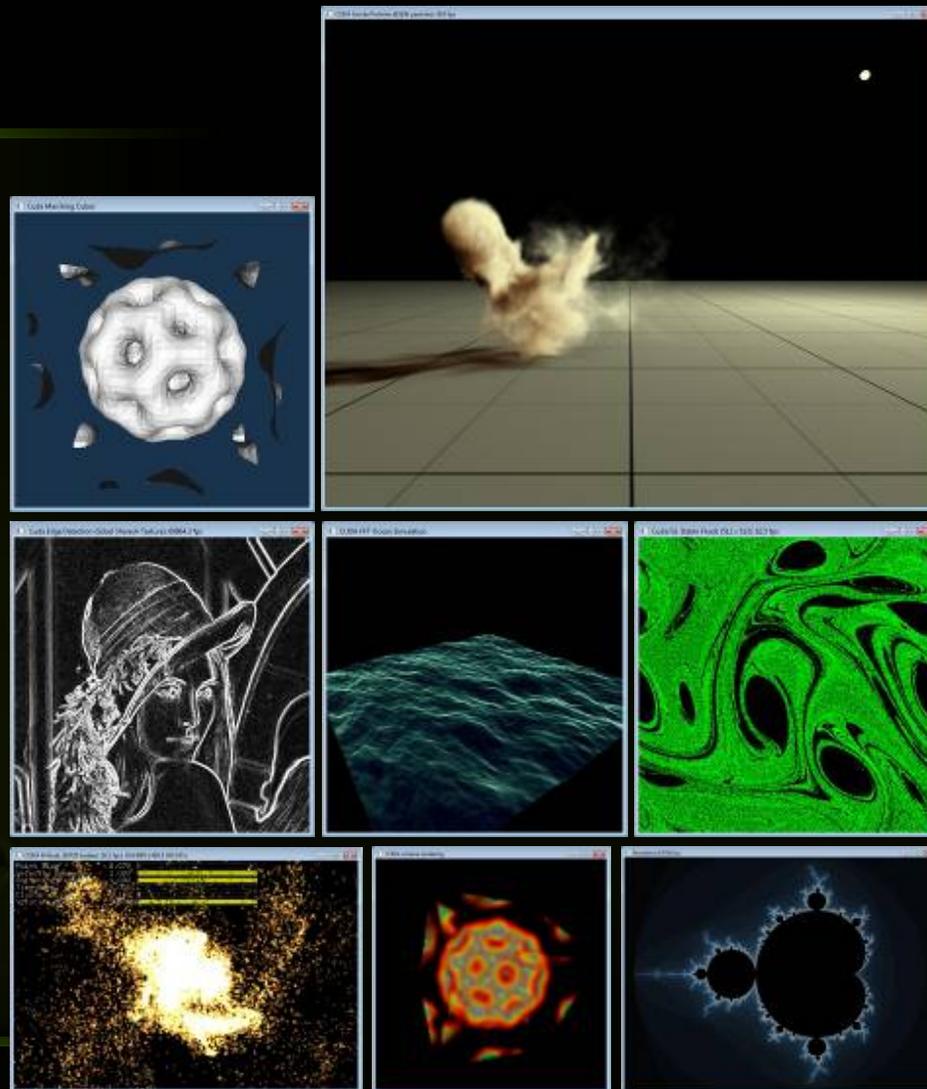
GPU computing and Graphics development best practices



# NVIDIA SDKs

Hundreds of code samples for  
CUDA C/C++, DirectCompute,  
and OpenCL

- Finance
- Oil & Gas
- Video/Image Processing
- 3D Volume Rendering
- Particle Simulations
- Fluid Simulations
- Math Functions



# NVIDIA Application Acceleration Engines

## OptiX – ray tracing engine

- Programmable GPU ray tracing pipeline that greatly accelerates general ray tracing tasks
- Supports programmable surfaces and custom ray data



*OptiX shader example*

## PhysX – physics and dynamics engine

- GPU computed physics – fast enough for real time.
- In use across games, DCC tools, and simulation



*PhysX particles example*

## SceniX – scene management engine

- High performance OpenGL scene graph built around CgFX for maximum interactive quality
- Provides ready access to new GPU capabilities & engines



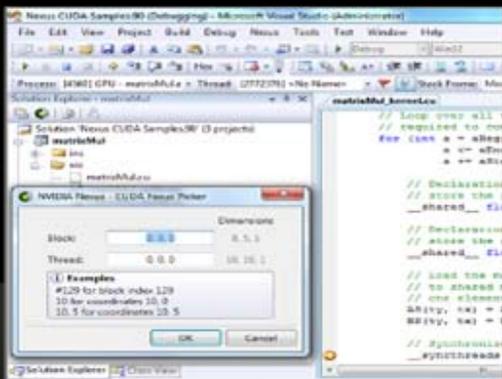
*Autodesk Showcase customer example*

# “Nexus” 1.0 Beta



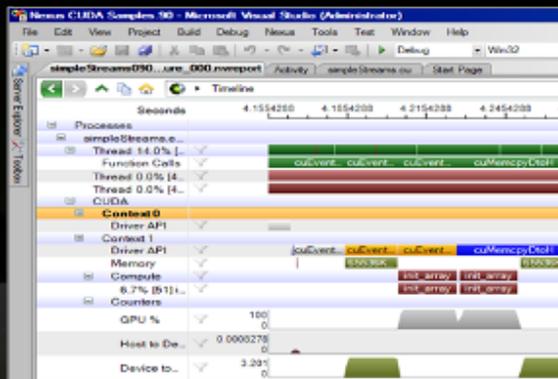
PARTNER

ПАРТНЕР



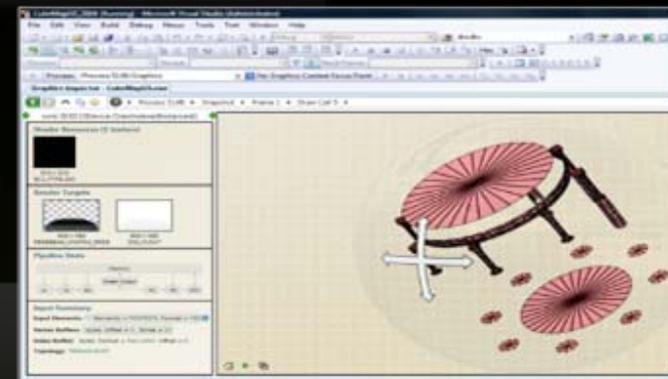
## Parallel Debugger

GPU source code debugging  
Variable & memory inspection



## System Analyzer

Platform-level Analysis  
For the CPU and GPU  
Visualize Compute Kernels,  
Driver API Calls, and  
Memory Transfers



## Graphics Inspector

Visualize and debug  
graphics content

# Massively Parallel Development

## Tools for Windows

	Visual Studio Integration	Parallel Debugging		Parallel Profiling		System Analysis/ Trace (CPU/GPU)		Premium Support	Early Access	Multi-Vendor GPU Support
		Compute	Gfx	Compute	Gfx	Compute	Gfx			
<b>“Nexus”* Pro</b>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<b>“Nexus”*</b>	✓	✓	✓	✓	✓					
<b>NV Visual Profiler</b>				✓						
<b>Mcore Platform Analyzer</b>				✓		✓				
<b>PerfKit/PerfHUD</b>			✓		✓		✓			
<b>gDEDebugger</b>			✓		✓		✓	✓		✓

\* “Nexus” is NVIDIA’s code name

# Massively Parallel Development

## Tools for Linux

	Compilation	Debugging	Profiling	Analysis	Premium Support
<b>CAPS HMPP</b>	✓				✓
<b>PGI Accelerators</b>	✓				✓
<b>PGI CUDA Fortran</b>	✓				✓
<b>NV cuda-gdb</b>		✓			
<b>Allinea DDT</b>		✓			✓
<b>TotalView Debugger</b>		✓			✓
<b>NV Visual Profiler</b>			✓		
<b>TAU CUDA</b>			✓	✓	
<b>Mcore Platform Analyzer</b>			✓	✓	✓

cuda-gdb will be extended to support OpenCL debugging in a future release, with solutions from Allinea, TotalView and others expected to follow.

# Massively Parallel Development

## Tools for Linux

	Debugging	Profiling	Analysis	Cluster Support	Lib's
CUDA C/C++	✓	✓	✓	✓	✓
OpenCL	Coming Soon	✓	Coming Soon	✓	Coming Soon
Fortran		✓	✓	✓	✓

cuda-gdb will be extended to support OpenCL debugging in a future release, with solutions from Allinea, TotalView and others expected to follow.

# cuda-gdb

## CUDA debugging **integrated** into GDB on Linux

- Supported on **32bit** and **64bit** systems
- **Seamlessly** debug both the host/CPU and device/GPU code
- Set **breakpoints** on any source line or symbol name
- Access and print all CUDA memory allocs, local, global, constant and shared vars

Included in the CUDA Toolkit

```
*gud-acos_dbg* - emacs@ssalian-linux
File Edit Options Buffers Tools Gud Complete In/Out Signals Help

p p* [debugger icons]

[Current CUDA Thread <<<(0,0),(>
acos_main () at /ssalian-local/
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at /ssalian-local/
Breakpoint 2 at 0x805abc4c: fil

#else /* FERMI */
acos_main<<<ACOS_CTA_CNT,ACOS_THREAD_CNT>>
#endif
stop = second();
cudaStat = cudaGetLastError(); /* check f

NX - ssalian@172.16.175.110:1022 - ssalian-linux
Applications Places System

*gud-acos_dbg* - emacs@ssalian-linux
File Edit Options Buffers Tools Gud Complete In/Out Signals Help

p p* [debugger icons]

[Current CUDA Thread <<<(0,0),(>
Breakpoint 1, acos_main () at a>
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at acos.cu:390
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at acos.cu:391
(cuda-gdb) p threadIdx
$5 = {x = 0, y = 0, z = 0}
(cuda-gdb) p blockIdx
$6 = {x = 0, y = 0}
(cuda-gdb) info cuda threads
<<<(0,0),(0,0,0)>> ... <<<(0,0)
<<<(0,0),(32,0,0)>> ... <<<(23)
(cuda-gdb) p blockDim
$7 = {x = 128, y = 1, z = 1}

__device_func__(float __cuda_acosf(float a))
{
float t0, t1, t2;

t0 = __cuda_fabsf(a);
t2 = 1.0f - t0;
t2 = 0.5f * t2;
t2 = __cuda_sqrtf(t2);
t1 = t0 > 0.57f ? t2 : t0;
t1 = __internal_asinf_kernel(t1);
t1 = t0 > 0.57f ? 2.0f * t1 : CUDART_PIO2_F
if (__cuda_signbit(a))
t1 = CUDART_PI2_F - t1;
}
#if !defined(__CUDABE__
if (__cuda_isnanf(a))
t1 = a + a;
}

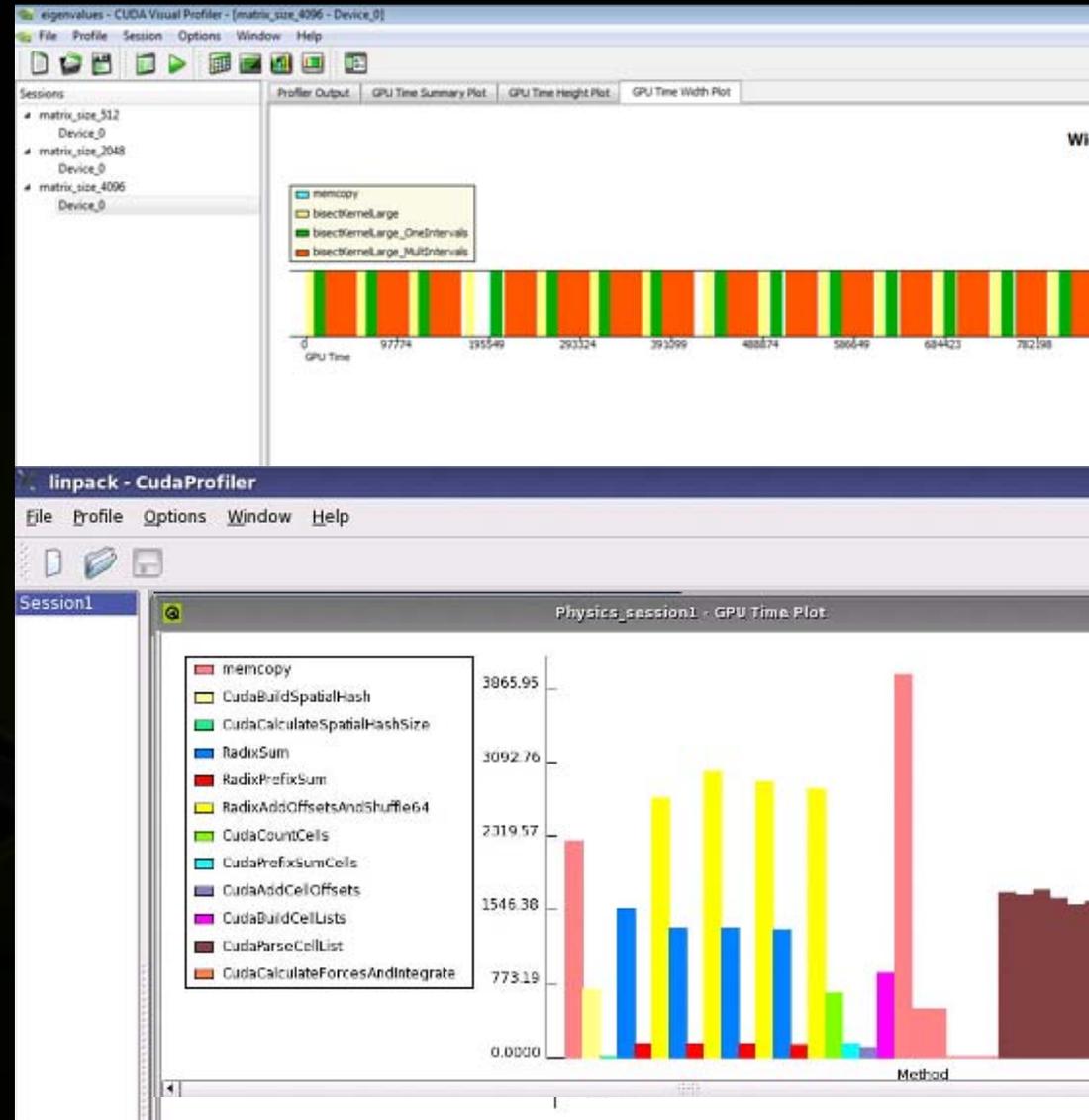
Kernel Source Debugging

$1 = {x = 128, y = 1, z = 1}
(cuda-gdb) b pthreadDim
<<<(0,0),(32,0,0)>> ... <<<(23)
(cuda-gdb) p blockDim
$7 = {x = 128, y = 1, z = 1}
(cuda-gdb) p threadIdx
$5 = {x = 0, y = 0, z = 0}
(cuda-gdb) p blockIdx
$6 = {x = 0, y = 0}
(cuda-gdb) info cuda threads
<<<(0,0),(0,0,0)>> ... <<<(0,0)
<<<(0,0),(32,0,0)>> ... <<<(23)
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at acos.cu:391
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at acos.cu:390
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at /ssalian-local/
(cuda-gdb) s
[Current CUDA Thread <<<(0,0),(>
acos_main () at /ssalian-local/
Breakpoint 2 at 0x805abc4c: fil
```

# CUDA Visual Profiler

- **Analyze GPU HW performance signals, kernel occupancy, instruction throughput, and more**
- **Highly configurable tables and graphical views**
- **Save/load profiler sessions or export to CSV for later analysis**
- **Compare results visually across multiple sessions to see improvements**
- **Windows, Linux and Mac OS X supported**  
OpenCL Visual Profiler for Windows and Linux

Included in the **CUDA Toolkit**



# Massively Parallel Development (summary)

## Tools for Linux - Language / API Support

	CUDA C/C++	OpenCL	Fortran
CAPS HMPP	✓		✓
PGI Accelerators	✓		✓
PGI CUDA Fortran	✓		✓
NV cuda-gdb	✓		
Allinea DDT	✓		
TotalView Debugger	✓		
NV Visual Profiler	✓	✓	
TAU CUDA	✓		✓
Mcore Platform Analyzer	✓		✓

cuda-gdb will be extended to support OpenCL debugging in a future release, with solutions from Allinea, TotalView and others expected to follow.

# 200+ Universities Teaching GPU Computing on CUDA



**Thank You**

