# NVIDIA Parallel Nsight™

**SuperComputing 2010 | New Orleans**

# Agenda

- **Introduction to Parallel Nsight**

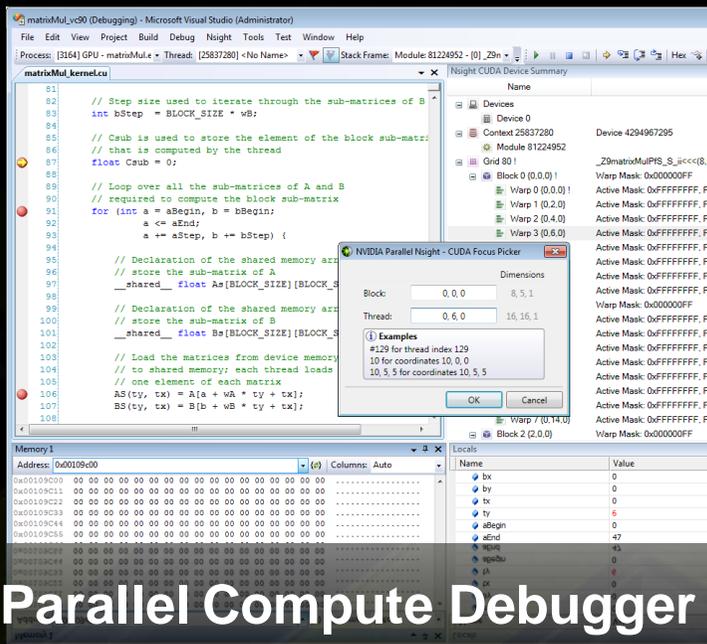- **CUDA C/C++ Source Debugging**

- **Analysis/System Trace**

- **Q&A**

# What is Parallel Nsight?

- Development environment for heterogeneous platforms (CPU and GPU)

- Fully integrated into Microsoft Visual Studio 2008 and 2010

- Dramatic productivity improvement in common development tasks

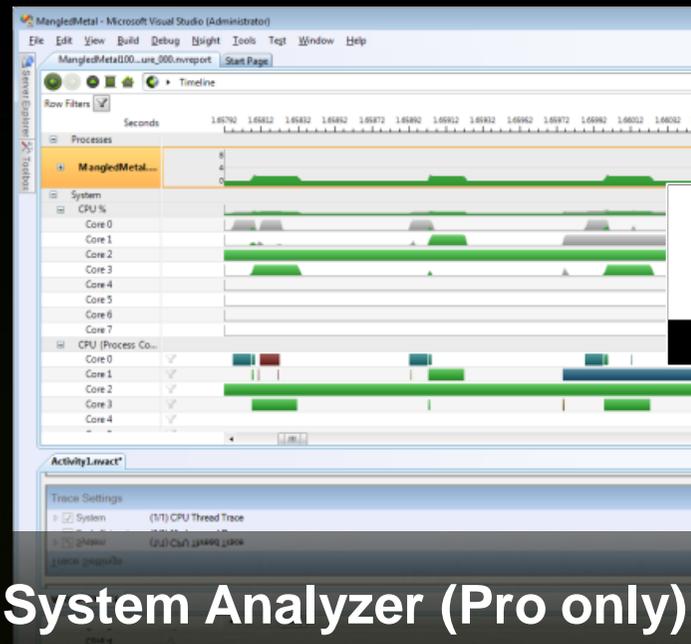# Parallel Nsight for Compute



## Parallel Compute Debugger

Examine compute kernels directly on GPU hardware

Debug CUDA C/C++ and DirectCompute applications

Visualize thousands of threads executing in parallel using Visual Studio

Use conditional breakpoints to correct errors in massively parallel code

## System Analyzer (Pro only)

Capture and visualize CPU and GPU level events on a single correlated timeline

Inspect workload dependencies using the Timeline View

Profile CUDA kernels using GPU performance counters

# CUDA C/C++ Debugging

- Compile your code with Debug flag

- Use the familiar Visual Studio interface to debug your GPU code

- Dramatic productivity improvements
  - Explore memory during a live session vs. coding specific transfers
  - Immediately view live variables vs. printf/recompile loop
  - Set data breakpoints on memory area vs. trial and error
  - And much more!

# Setting Breakpoints

# Viewing Variable Values

| Name | Value | Type |
|------|-------|------|
| a + wA * ty + tx | 16 | int |
| C[a+ wA *ty + tx] | 13.683188 | __device__ float& |
| C[a + wA * ty + tx + 10] | 13.49013 | __device__ float& |
| C[16] + 34 | 47.683189 | float |

Watch 1

Autos    Locals    Watch 1

# Viewing GPU Memory

# Switching Between Threads

# Conditional Breakpoints

# Data Breakpoints

# Trace

- **Powerful performance tool**

- **Correlated timeline between CPU and GPU**

- **Bottleneck identification**
  - **Find CPU vs. GPU boundedness**
  - **Find memory transfer vs. kernel computation bounded**
  - **Get macro-level information on which CUDA kernels use the most time**

# Trace across the CPU and GPU

# Trace – Overlap Memory Transactions

# Tesla Compute Cluster Support (TCC)

- TCC is a special driver mode for Windows 7, Vista, and HPC Server 2008.
  - Included in our most recent R260 driver release.

- Benefits
  - Lower overhead to kernel launches = higher performance
  - Running CUDA on a MS Remote desktop

- Parallel Nsight 1.5 now supports debugging on GPUs using a TCC driver.

# 4 Flexible GPU Development Configurations

**Desktop**

**Single machine, Single NVIDIA GPU**

Analyzer, Graphics Inspector

**NEW** **Single machine, Dual NVIDIA GPUs**

Analyzer, Graphics Inspector, Compute Debugger

**Networked**

**Two machines connected over the network**

Analyzer, Graphics Inspector, Compute Debugger, Graphics Debugger

TCP/IP

**Workstation SLI**

**SLI Multi OS workstation with two Quadro GPUs**

Analyzer, Graphics Inspector, Compute Debugger, Graphics Debugger

# Parallel Nsight 1.5 Feature Support

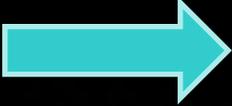| | Standard (no cost) | Professional ($349) *available for purchase in December* |
|---|:---:|:---:|
| Compute Debugger | ✓ | ✓ |
| DirectX 10 & 11 Debugger & Graphics inspector | ✓ | ✓ |
| GeForce Support:  9 series or higher | ✓ | ✓ |
| Tesla Support: C1050/S1070 or higher | ✓ | ✓ |
| Quadro Support: G9x or higher | ✓ | ✓ |
| Windows 7, Vista and HPC Server 2008 | ✓ | ✓ |
| Visual Studio 2008 SP1 and Visual Studio 2010 | ✓ | ✓ |
| Compute Analyzer | | ✓ |
| OpenGL and OpenCL Analyzer | | ✓ |
| DirectX 10 & 11 Analyzer | | ✓ |
| Tesla Compute Cluster (TCC) Debugging | | ✓ |

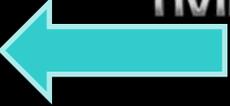http://www.nvidia.com/GetParalleINsight

# Parallel Nsight Resources

- **Parallel Nsight GPU Computing Forum**

- **The Parallel Nsight User Guide**
  - Installed with the Host installer
  - Available on the Web

- **Links to these from: http://developer.nvidia.com/ParalleINsight**