# NVIDIA GTX200: TeraFLOPS Visual Computing

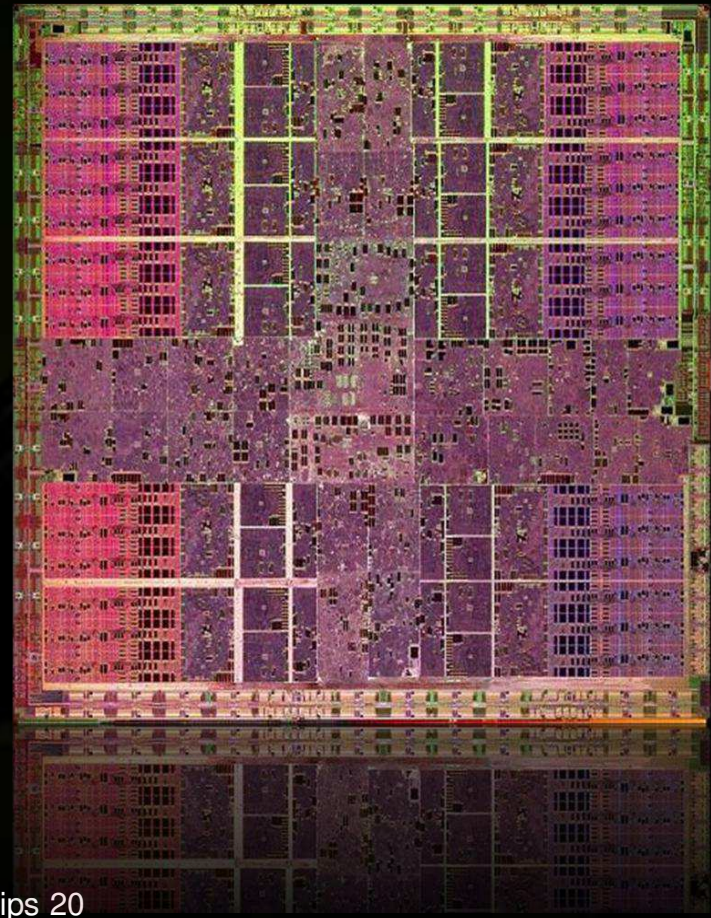**August 26, 2008**
**John Tynefield**
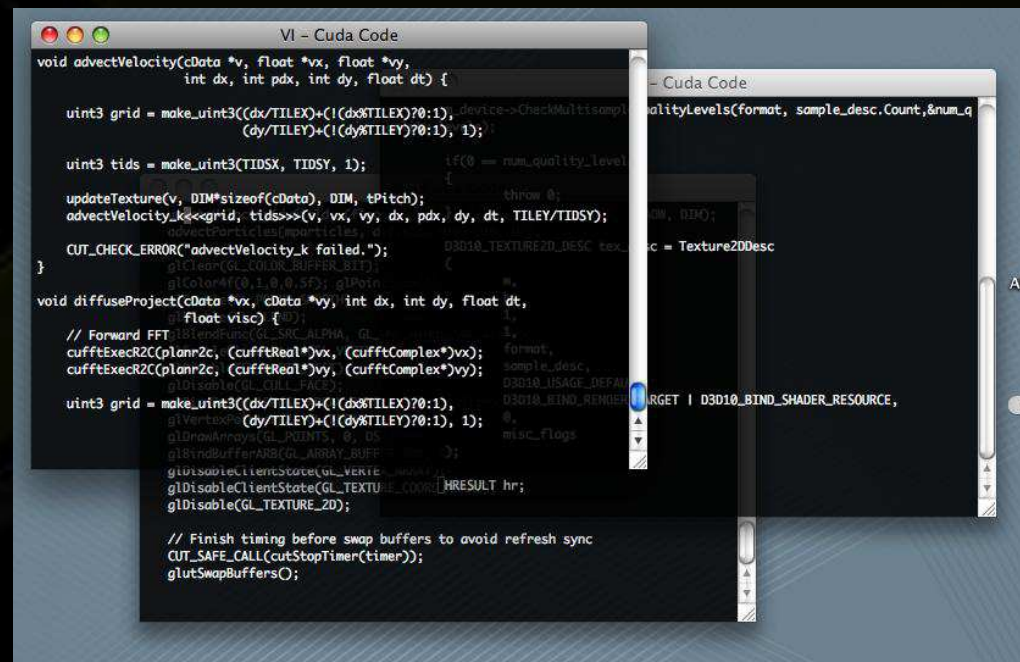
# Outline

- **Execution Model**
- **Architecture**
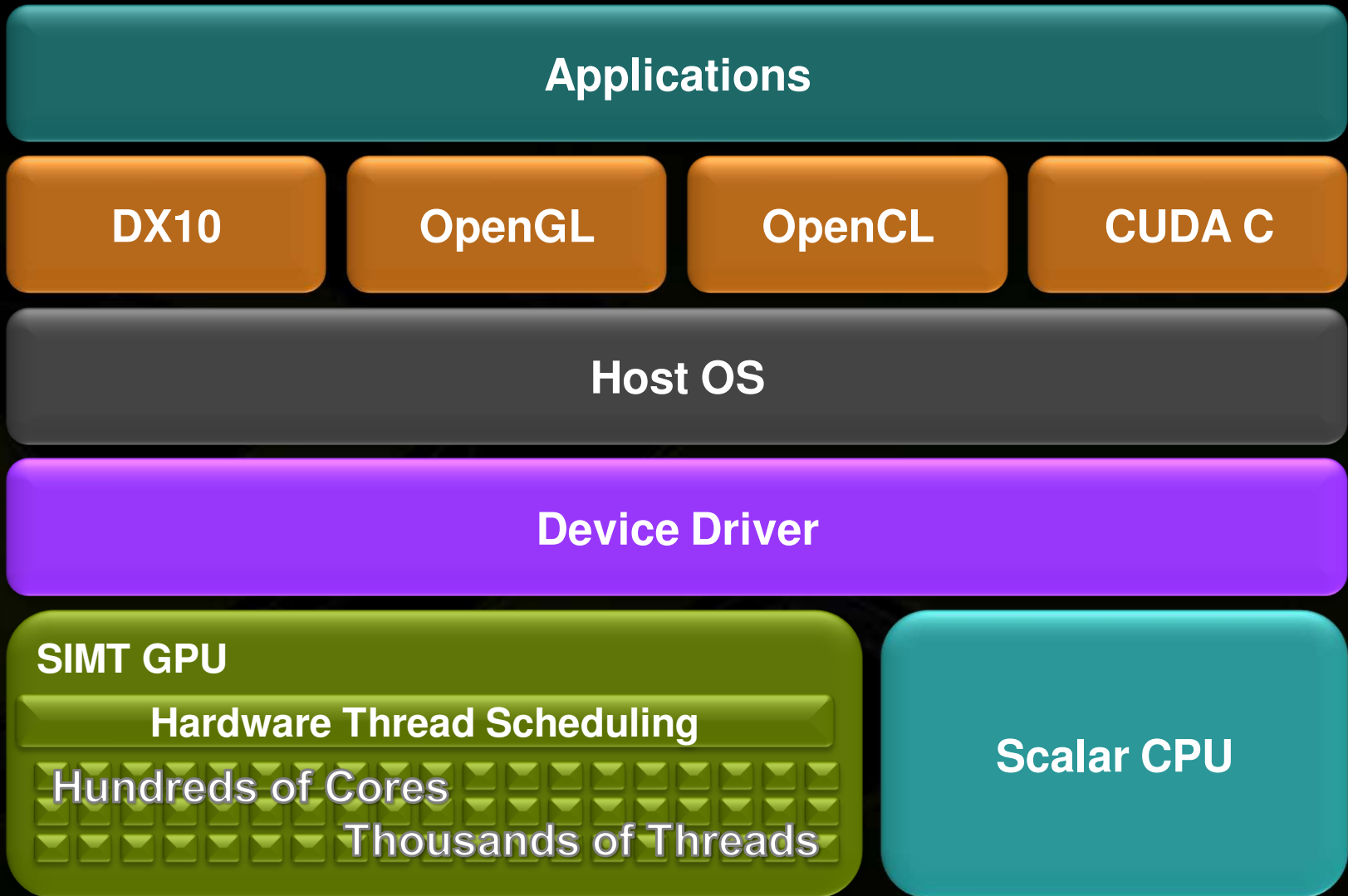- **Demo**

GTX200 – Hot Chips 20

# *Execution Model*

GTX200 – Hot Chips 20

# Software Architecture

**Applications**

**DX10**     **OpenGL**     **OpenCL**     **CUDA C**

**Host OS**

**Device Driver**

**SIMT GPU**

**Hardware Thread Scheduling**

Hundreds of Cores

Thousands of Threads

**Scalar CPU**

# SIMT Multithreaded Execution

**Single-Instruction Multi-Thread instruction scheduler**

time

| warp 8 instruction 11 |
| warp 1 instruction 42 |
| warp 3 instruction 95 |
| warp 8 instruction 12 |
| warp 3 instruction 96 |

- **SIMT: Single-Instruction Multi-Thread applies instruction to independent threads**

- **SIMT provides easy single-thread scalar programming with SIMD efficiency**

- **Warp: the set of 32 parallel threads that execute a SIMT instruction**

- **Hardware implements zero-overhead warp and thread scheduling**

- **Each thread processor manages state for up to 128 threads**

- **SIMT threads execute independently**

- **SIMT warp diverges and converges when threads branch independently**

- **Best efficiency and performance when threads of a warp execute together**

# Execution Model - CUDA

**Thread**

Local Memory

Register File

**Block**

Shared Memory

- Local Memory:   per-thread
  - Private per thread
  - Auto variables, register spill
- Shared Memory:   per-block
  - Shared by threads of CTA
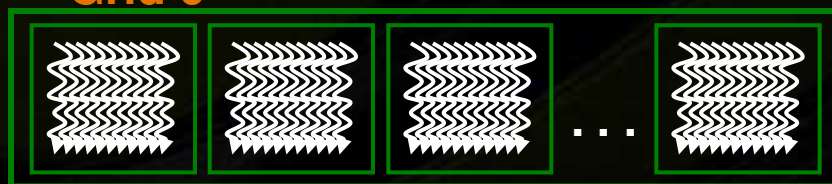  - Inter-thread communication
- Global Memory:   per-application
  - Shared by all threads
  - Inter-Grid communication

**Grid 0**

. . .

**Grid 1**

. . .

Global Memory

Sequential Grids in Time

# Execution Model - Graphics

**Pixel Thread**

**Register File**

**Pixel Warp**

**Plane Equations**

**Texture**

- Local Registers: per Pixel
  - Private per pixel
- Planes/Textures: per Warp
  - Define surface-space inputs
  - Array of 1D/2D/3D data arrays
- Target Images: per Grid
  - Array of 2D surfaces

**Pixel Grid 0**

. . .

**Target Images**

**Pixel Grid 1**

. . .

**Target Images**

# *Architecture*

GTX200 – Hot Chips 20

# Implementation

- **TSMC 65nm**
- **1.4 Billion Transistors**
- **24.3 x 24.0 mm die**
- **2236 Ball BGA**

- **1 TeraFLOPS SP / 84 GigaFLOPS DP**
  - **1.4 GHz Processor Clock**
- **140 GB/s**
  - **1.1 GHz Memory Clock**
- **Up to 4GB on-board Memory**

GTX200 – Hot Chips 20

# GTX200 Unified Visual Computing

**Processing Cores**

**Shared Memory**

**Processor Communication Fabric**

Fixed Function Acceleration

Memory & I/O

- Tesla Unified Graphics and Computing Architecture
- Scales parallel performance 2X beyond G80
- 240 ▪ Thread Processor cores, 30K threads
- Double precision 64-bit IEEE 754 floating point

# SM Multithreaded Multiprocessor

**I - Cache**

**MT Issue**

**C - Cache**

| SP | SP |
| SP | SP |
| SP | SP |
| SP | SP |
| SFU | SFU |

**DP**

**Shared Memory**

- **8 SP Thread Processors**
  - IEEE 754 32-bit floating point
  - 32-bit and 64-bit integer
  - 2K 32-bit registers per SP
  - Variable 4-128 registers / thread
- **2 SFU Special Function Units**
- **1 DP Double Precision Unit**
  - IEEE 754R 64-bit floating point
  - Fused multiply-add
- **Scalar register-based ISA**
- **Multithreaded Instruction Unit**
  - 1024 threads, hardware multithreaded
  - 32 SIMT warps of 32 threads
  - Independent thread execution
  - Mixed concurrent thread types
- **16KB Shared Memory**
  - Concurrent threads share data
  - Low latency load/store

GTX200 – Hot Chips 20

# Double Precision Fused Multiply-Add

- **Revised IEEE 754R standard specifies FMA**
- **FMA has lower latency than FMAD**
- **FMA improves accuracy of many computations**
  - **Dot products**
  - **Matrix multiplication**
- **Enables fast mixed-precision convergence algorithms**
- **Enables higher performance algorithms for extended-precision arithmetic**
- **Enables efficient implementation of exactly-rounded division and square root**

| I - Cache |
|:---:|
| MT Issue |
| C - Cache |

| SP | SP |
|:---:|:---:|
| SP | SP |
| SP | SP |
| SP | SP |
| SFU | SFU |

| DP |
|:---:|

| Shared Memory |
|:---:|

# Coalesced Memory I/O

- When 16 threads access a contiguous region of device memory
- 16 data elements loaded in one instruction
  - int, float: 64 bytes (fastest)
  - int2, float2: 128 bytes
  - int4, float4: 256 bytes (2 transactions)
- Regions aligned to multiple of size
- Coalescing scales gracefully for partially contiguous accesses

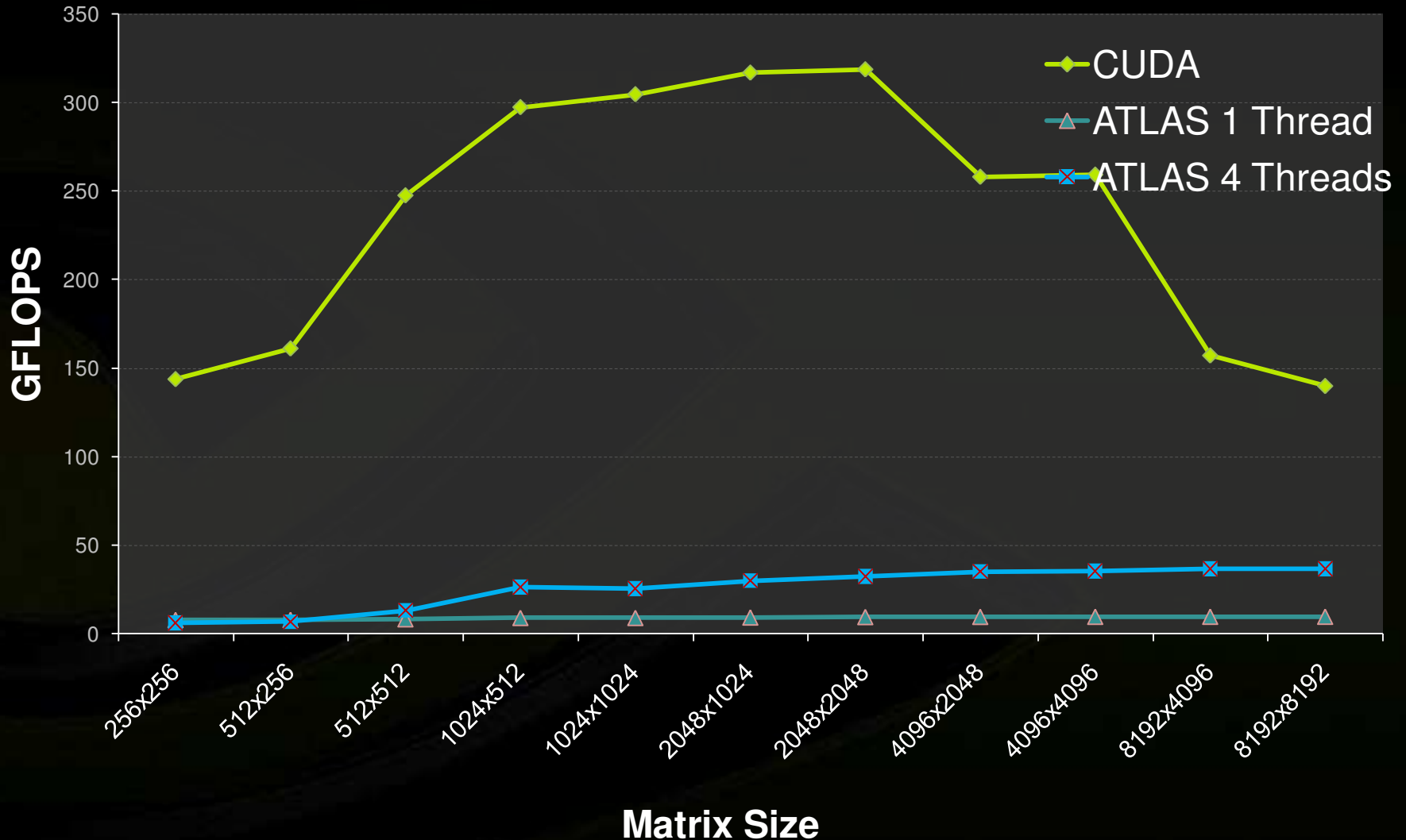| | |
|---|---|
| | Address 120 |
| | Address 124 |
| Thread 0 | Address 128 |
| Thread 1 | Address 132 |
| Thread 2 | Address 136 |
| Thread 3 | Address 140 |
| Thread 4 | Address 144 |
| Thread 5 | Address 148 |
| Thread 6 | Address 152 |
| Thread 7 | Address 156 |
| Thread 8 | Address 160 |
| Thread 9 | Address 164 |
| Thread 10 | Address 168 |
| Thread 11 | Address 172 |
| Thread 12 | Address 176 |
| Thread 13 | Address 180 |
| Thread 14 | Address 184 |
| Thread 15 | Address 188 |
| | Address 192 |
| | Address 196 |

64B segment

GTX200 – Hot Chips 20

# Performance: Single Precision BLAS

**BLAS (SGEMM) on CUDA**

CUBLAS: CUDA 2.0b2, Tesla C1060
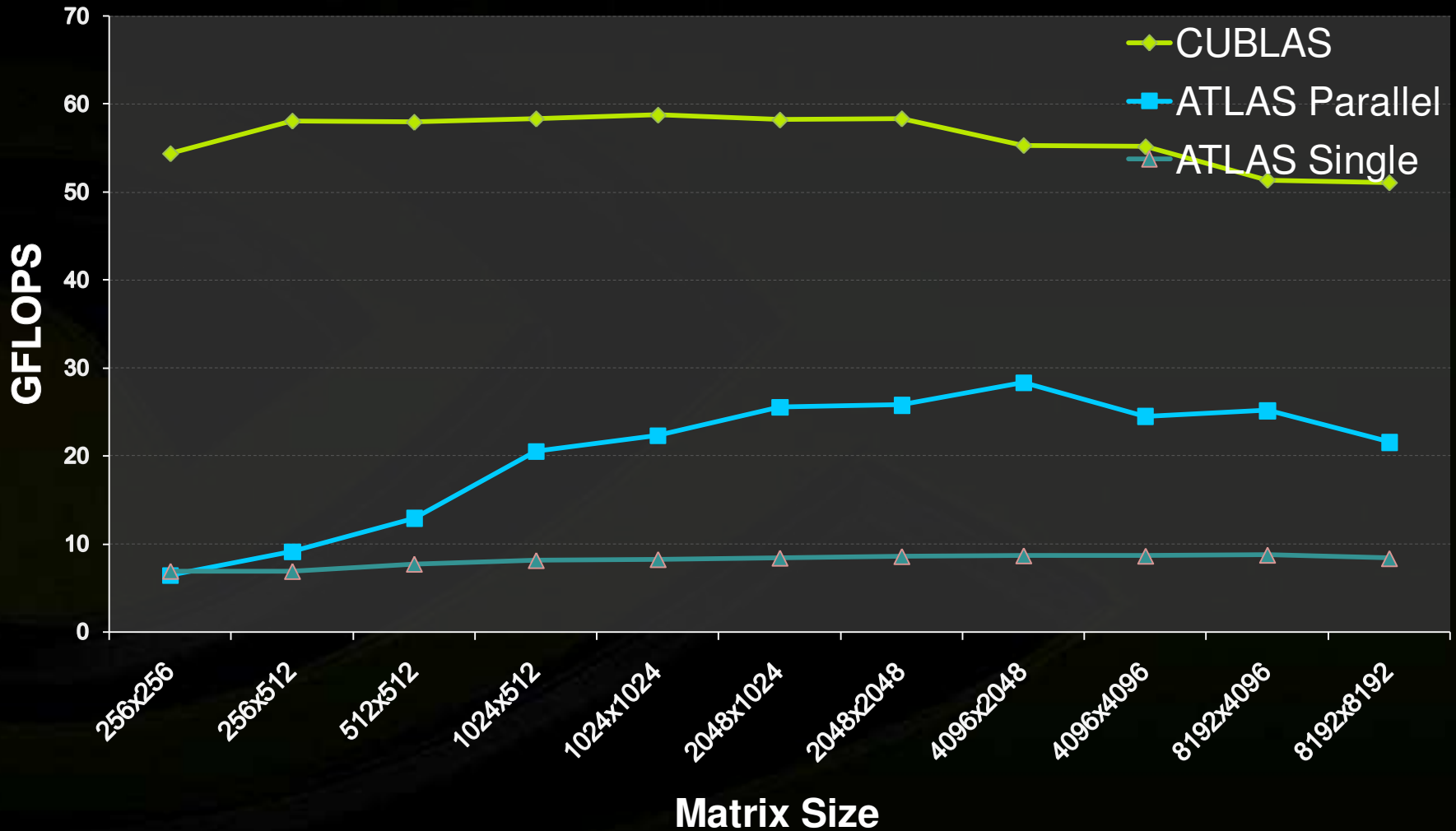ATLAS 3.81 on Dual 2.8GHz Opteron Dual-Core

GTX200 – Hot Chips 20

# Performance: Double Precision BLAS

**BLAS (DGEMM) on CUDA**

CUBLAS CUDA 2.0b2 on Tesla C1060
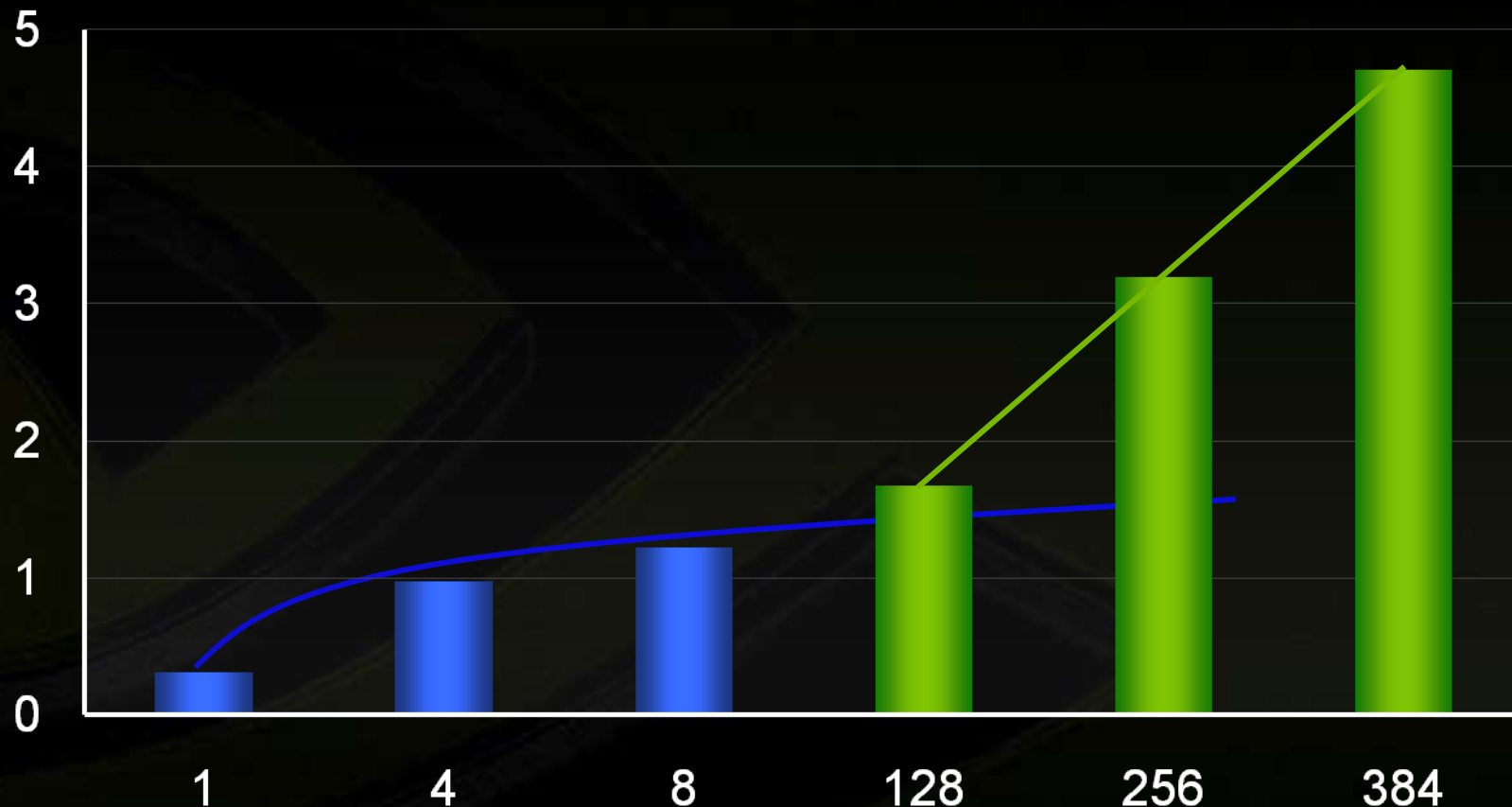ATLAS 3.81 on Intel Xeon E5440 Quad-core, 2.83 GHz



Legend:
- CUBLAS
- ATLAS Parallel
- ATLAS Single

Y-axis: **GFLOPS** (0–70)
X-axis: **Matrix Size** (256x256, 256x512, 512x512, 1024x512, 1024x1024, 2048x1024, 2048x2048, 4096x2048, 4096x4096, 8192x4096, 8192x8192)

GTX200 – Hot Chips 20

# Performance: Scaling

## Oil and Gas Computing: Reverse Time Migration
## Hand Optimized SSE Versus CUDA C
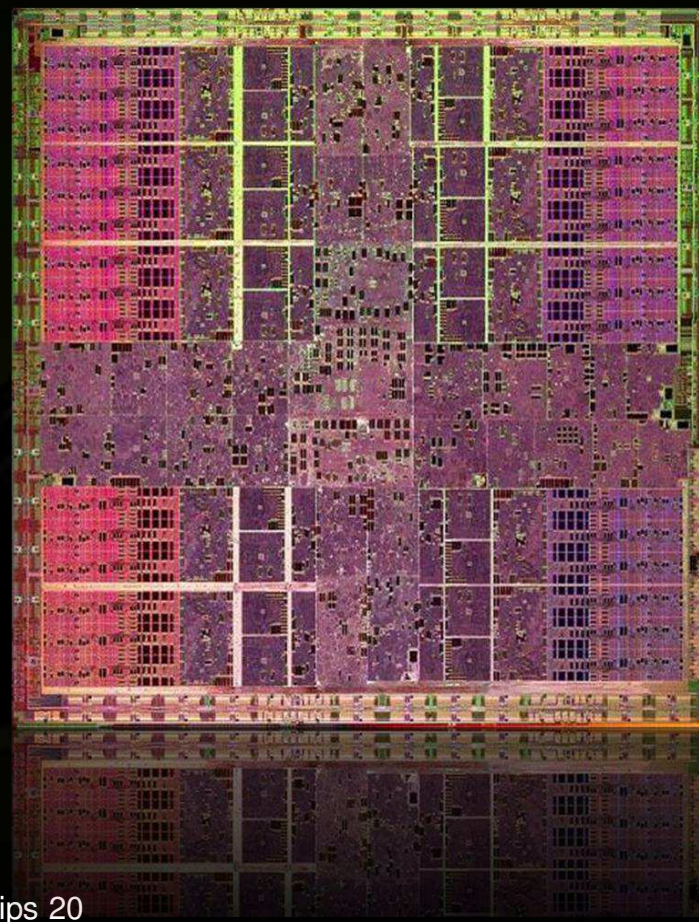


Billions of Points / sec

Number of Cores

X86 CPU

NVIDIA GPU

GTX200 – Hot Chips 20

# Summary

- **Rebalanced architecture to workload trends**
- **Scaled from 128 to 240 processors**
- **Hardware manages thousands of threads**
    - **Zero software overhead**
    - **Hides huge latencies**
    - **High achieved utilization**
- **Natively Scalar**
    - **No swizzling or vectorization overhead**
    - **Coalescing for high bandwidth memory I/O**
- **Software architecture allows 2X scaling on customer C code with no modification**

*Demo*

GTX200 – Hot Chips 20

# More Information

**Erik Lindholm, et. al., "NVIDIA Tesla: A Unified Graphics and Computing Architecture,"** *IEEE Micro***, March-April 2008.**

**http://www.nvidia.com/cuda**