# Power Consumption of Desktop and Mobile GPUs for IRSTAP Applications

Michael Roeder, Jeremy Furtek, Nolan Davis,
Cezario Tebcherani, Masatoshi Tanida and Dennis Braunreiter

High Performance Embedded Computing (HPEC) Workshop, 23-25 September 2008
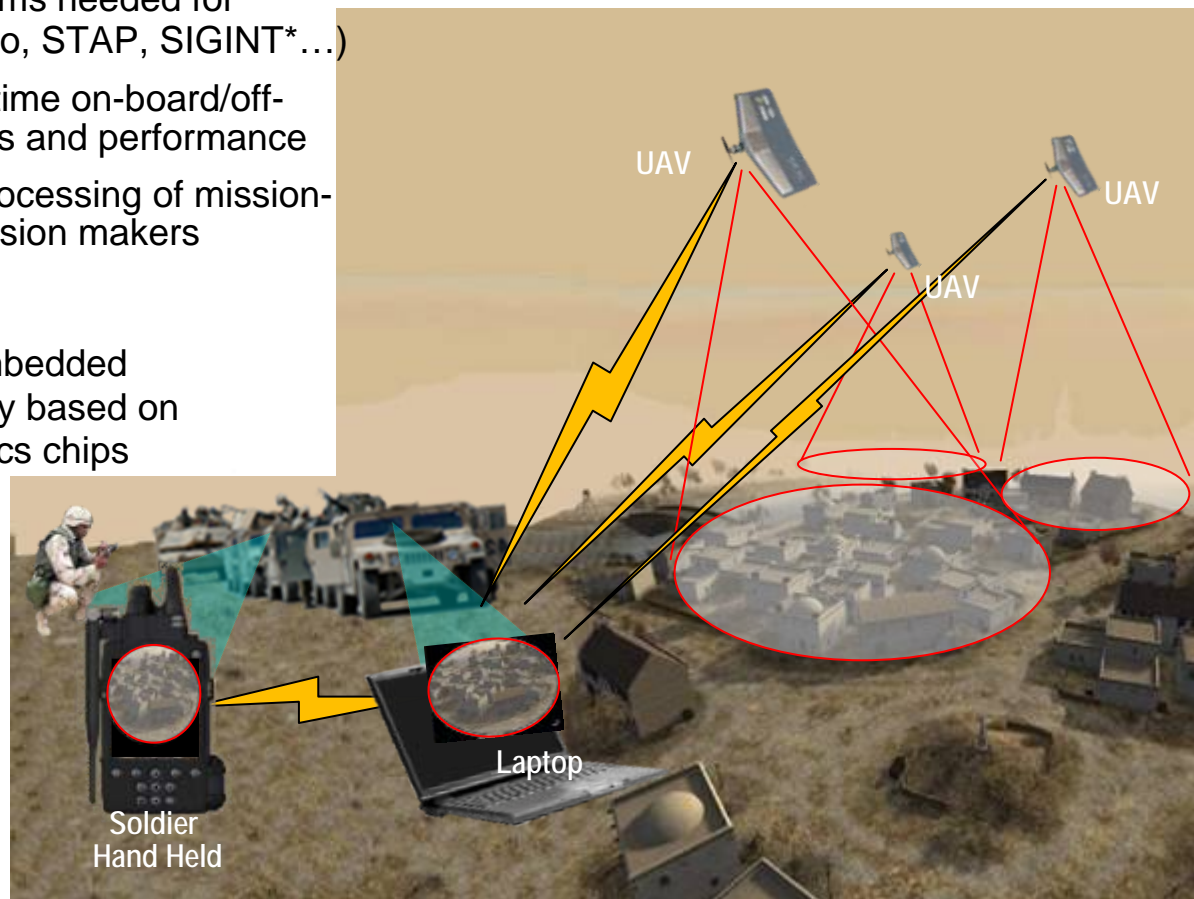
# STAP-BOY: Concept

**Problem:**

- Complex sensor modalities and algorithms needed for smaller platforms (SAR, 3D-motion video, STAP, SIGINT*…)

- Low-cost platform constraints limit real-time on-board/off-board and distributed sensing algorithms and performance

- Timely distribution, visualization, and processing of mission-critical data not available to tactical decision makers

**STAP-BOY goal:**

- Develop low-cost, scalable, teraflop, embedded multi-modal sensor processing capability based on commercial off-the-shelf (COTS) graphics chips
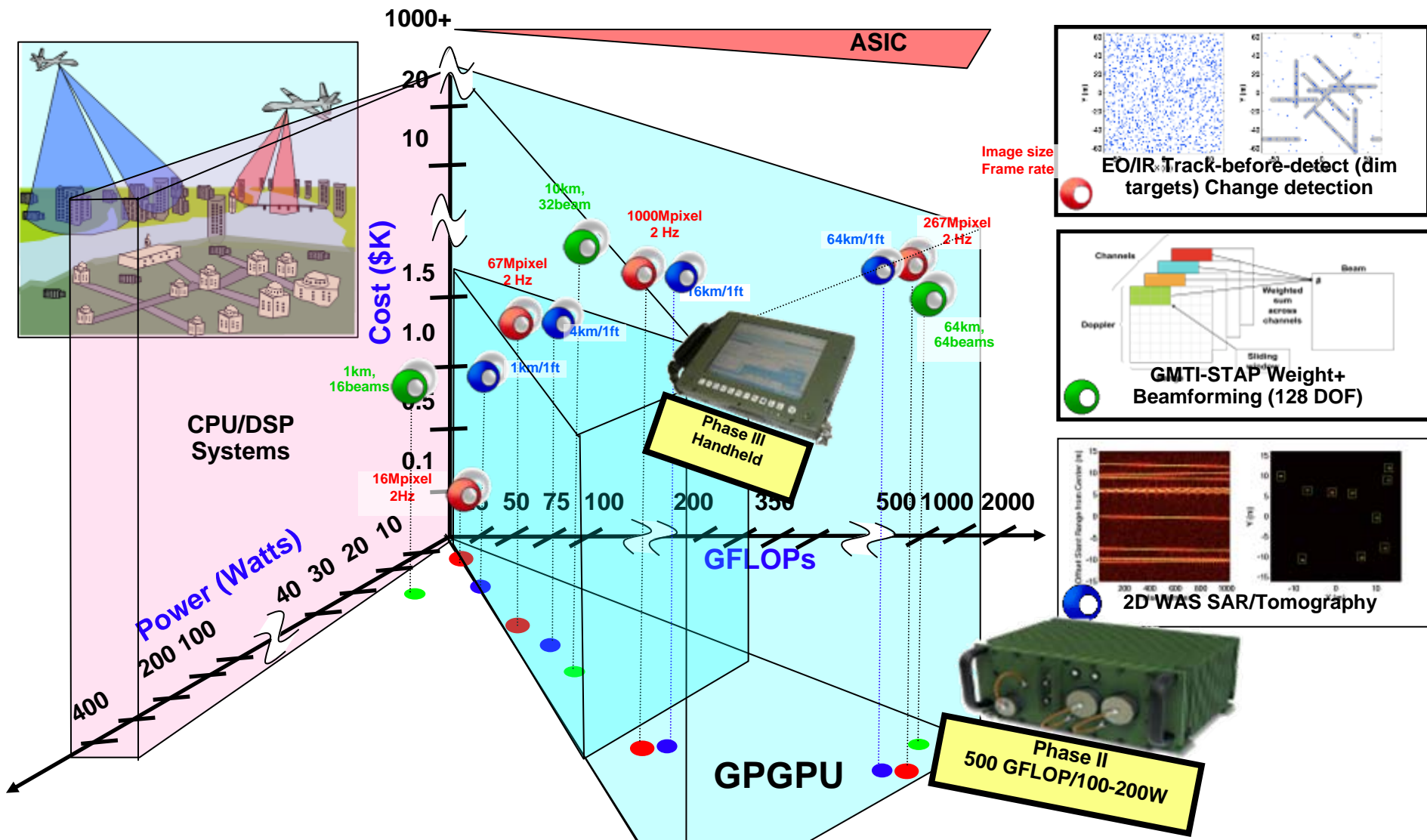
**STAP-BOY approach:**

- Map complex algorithms to COTS graphics chips with open source graphics languages

- Prototype scalable parallel embedded computing architecture for handhelds to teraflop single card

- Demonstrate on available tactically representative sensor systems

UAV

UAV

UAV

Laptop

Soldier
Hand Held

**Goal:**
½ Teraflop
10 Mobile GPUs
100W Total Power
$<15K

2

# Applications Pull



ASIC

CPU/DSP Systems

Cost ($K)

1000+
20
10
1.5
1.0
0.5
0.1

10km, 32beam

1000Mpixel 2 Hz

267Mpixel 2 Hz

64km/1ft

67Mpixel 2 Hz

16km/1ft

4km/1ft

64km, 64beams

1km, 16beams

1km/1ft

Phase III Handheld

16Mpixel 2Hz

Power (Watts)

400    200  100    40  30  20  10

50  75  100    200    350    500 1000 2000

GFLOPs

GPGPU

Phase II
500 GFLOP/100-200W

Image size
Frame rate

EO/IR Track-before-detect (dim targets) Change detection

Channels    Beam
Doppler    Weighted sum across channels    Sliding window

GMTI-STAP Weight+ Beamforming (128 DOF)

2D WAS SAR/Tomography

**Increasing Advanced Sensor Resolution Driving Processing Throughput Beyond CPU/DSP Capability In Communications and Cost-constrained Platforms**

3

SAIC
From Science to Solutions

# NVIDIA® GPU Performance/Power Study

| | GeForce™ 280 GTX | GeForce™ 8800 GTX | GeForce™ 8800M GTX |
|---|---|---|---|
| Category | Desktop | Desktop | Mobile |
| Process (nm) | 65 | 90 | 65 |
| Transistor Count (Millions) | 1400 | 681 | 754 |
| Stream Processors | 240 | 128 | 96 |
| Shader Clock (GHz) | 1.30 | 1.35 | 1.25 |
| Memory Clock (MHz) | 1107 | 900 | 800 |
| Memory Amount (MB) | 1024 | 768 | 512 |
| Thermal Design Power (W) | 236 | 177 | 65 |
| Release Date | June 2008 | November 2006 | November 2007 |

NVIDIA and GeForce are trademarks or registered trademarks of NVIDIA Corporation
in the United States and/or other countries.

SAIC
From Science to Solutions

# CPUs vs. GPUs: Head to Head (Floating Point)

Intel® quad-core Q9650

NVIDIA® 280 GTX

| Intel® quad-core Q9650 | | NVIDIA® 280 GTX |
|---|---|---|
| 820 million | **Transistors** | 1400 million |
| 3.0 GHz | **Clock Speed** | 1.3 GHz |
| 4 | **Number of Cores** | 240 |
| Serial | **Programming Model** | Highly parallel |
| Minimize latency | **Design Goal** | Maximize throughput |
| Complex cores:<br>• Branch prediction<br>• Out-of-order execution | **Design Approach** | Simple cores:<br>• Smaller caches<br>• In-order execution |
| 130 W | **Thermal Design Power (TDP)** | 236 W |
| 96 GFLOPS | **Theoretical Max. Computation Rate (single precision)** | 933 GFLOPS |

Intel is a registered trademark of Intel Corporation in the United States and/or other countries.
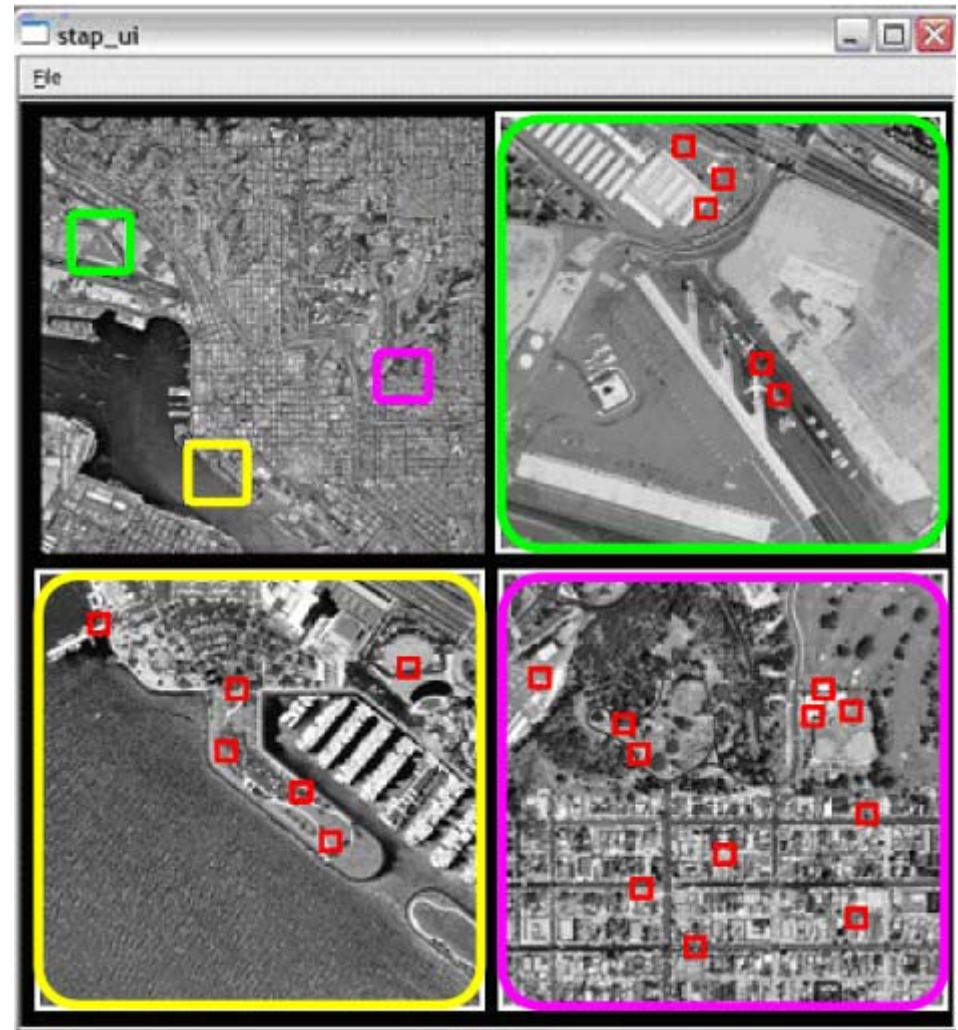NVIDIA is a registered is a registered trademark of NVIDIA Corporation in the United States and/or other countries.
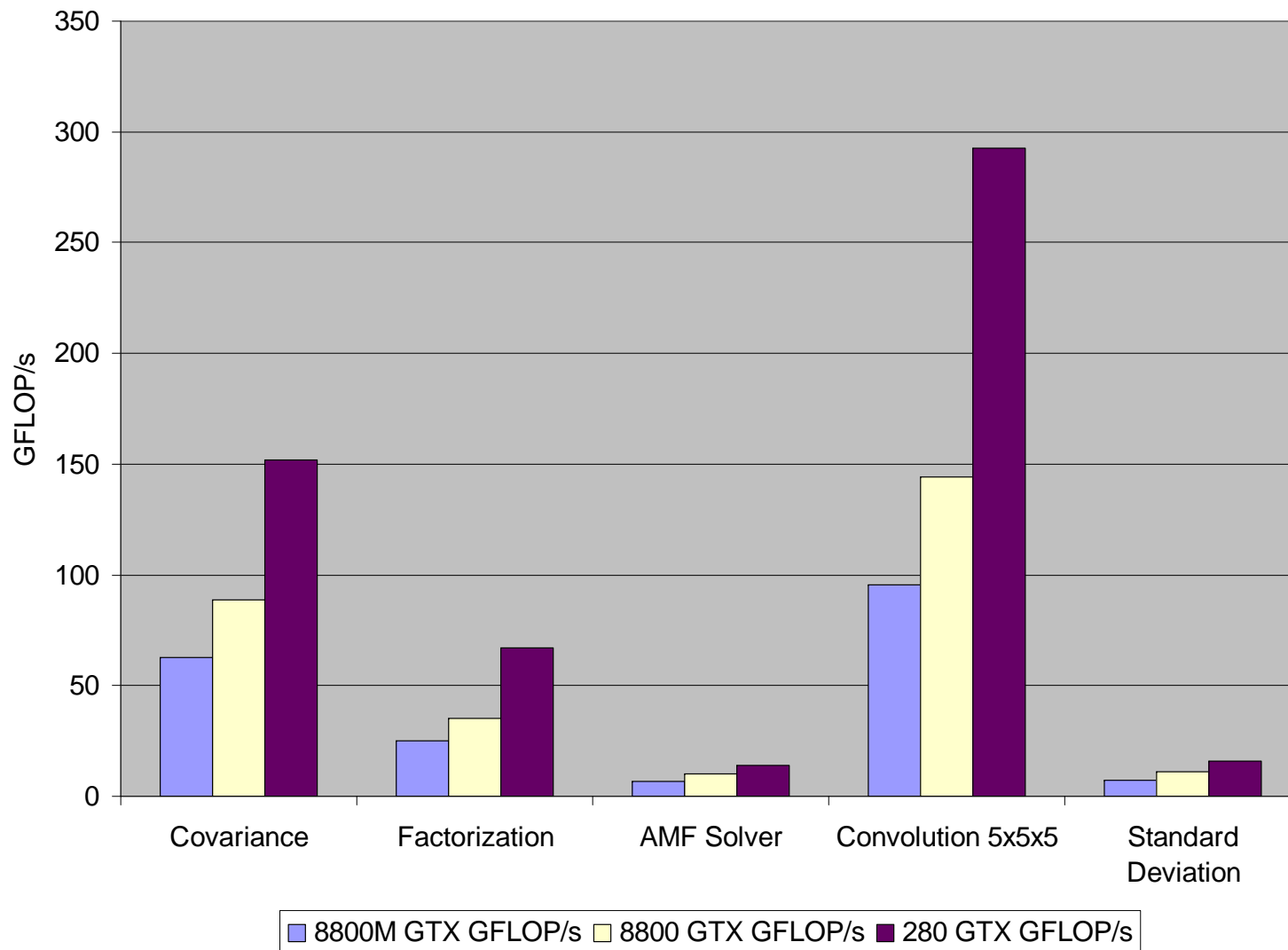
SAIC
*From Science to Solutions*

# IRSTAP Application

**Adapt Matched Filter Weights**

**STAP Velocity Hypothesis**

| Distribute Imagery to GPUs | Estimate Local Space-Time Covariance | Space-Time AMF | Compute Adaptive CFAR Thresholds | Gather Detections |
|---|---|---|---|---|

**Covariance Factorization**

**Space-Time Multiple Weight Solver**

**Space-Time AMF**

**Space-Time AMF**

**Apply Adaptive Matched Filters (AMF)**

**Tracker**

- **Adapt weights**
  - Covariance estimation
  - Covariance factorization
  - Find LMS solution
- **Apply AMF - convolution**
- **Compute CFAR thresholds**
  - Estimate standard deviation for each block
  - Adapt local standard deviations, excluding CUT
- **Gather detections into dense detection list**

6

*SAIC*
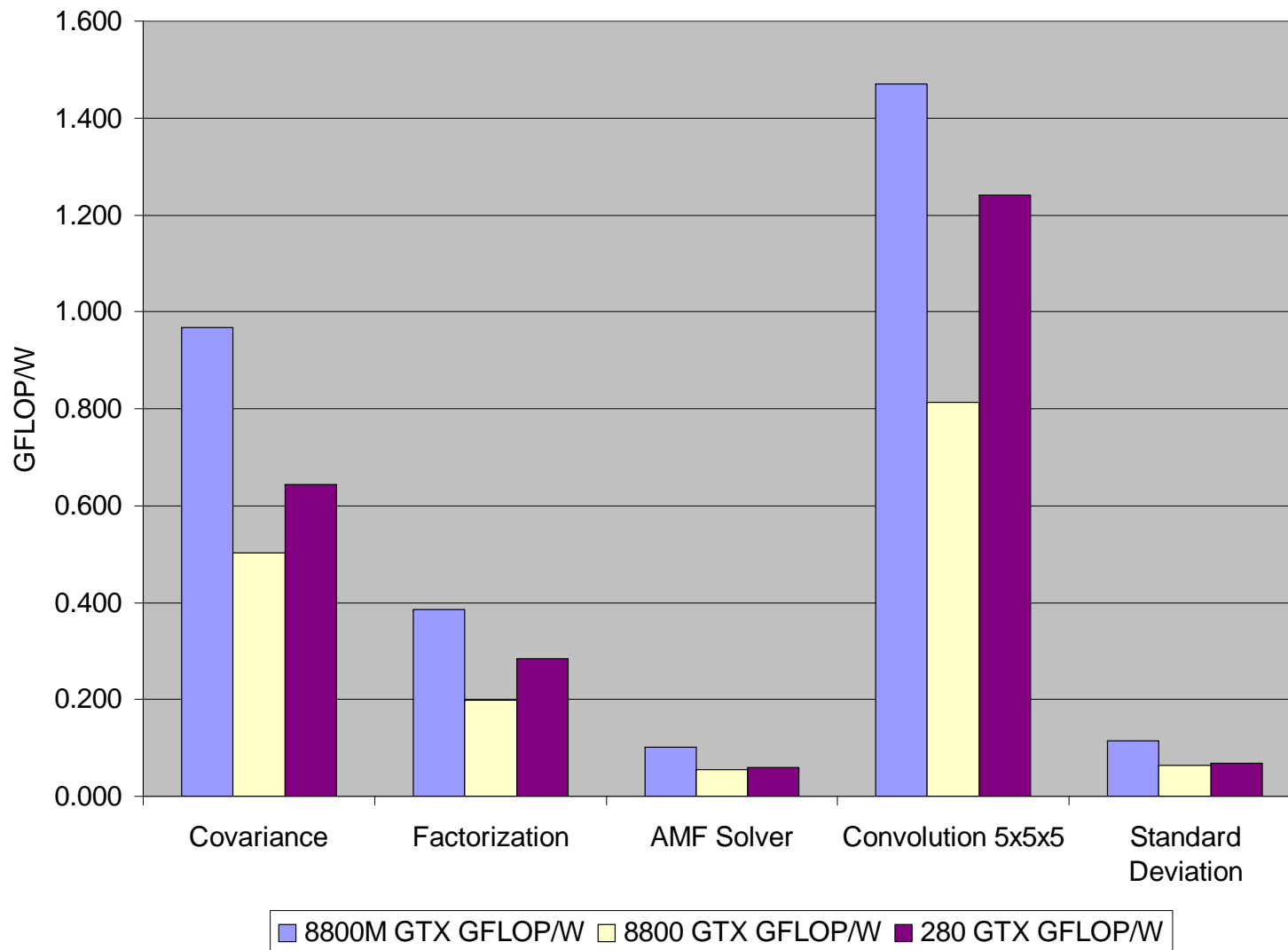*From Science to Solutions*

# IRSTAP Kernels

- **Image preprocessing**
  - Registration
  - Clutter subtraction
  - Spatial demeaning
- **Covariance calculation**
  - Ensemble generation
  - Covariance calculation
- **Weight adaptation**
  - Cholesky Factorization
  - Back and forward substitution
- **Weight application**
  - Convolution
- **Adaptive CFAR threshold adaptation**
  - Localized standard deviation
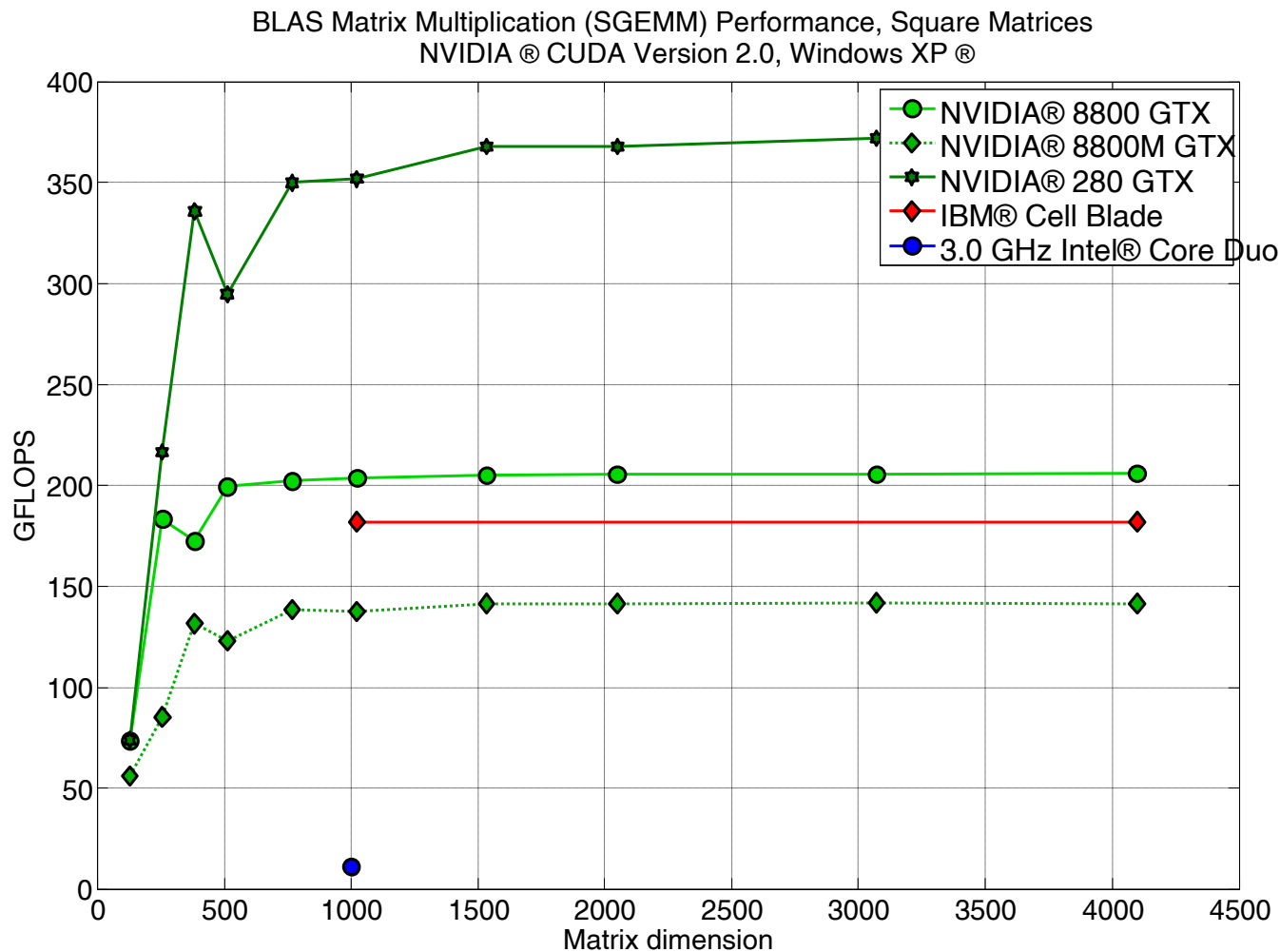  - Complex mean kernel
- **Detection reduction**

# IRSTAP Benchmarks GFLOP/s
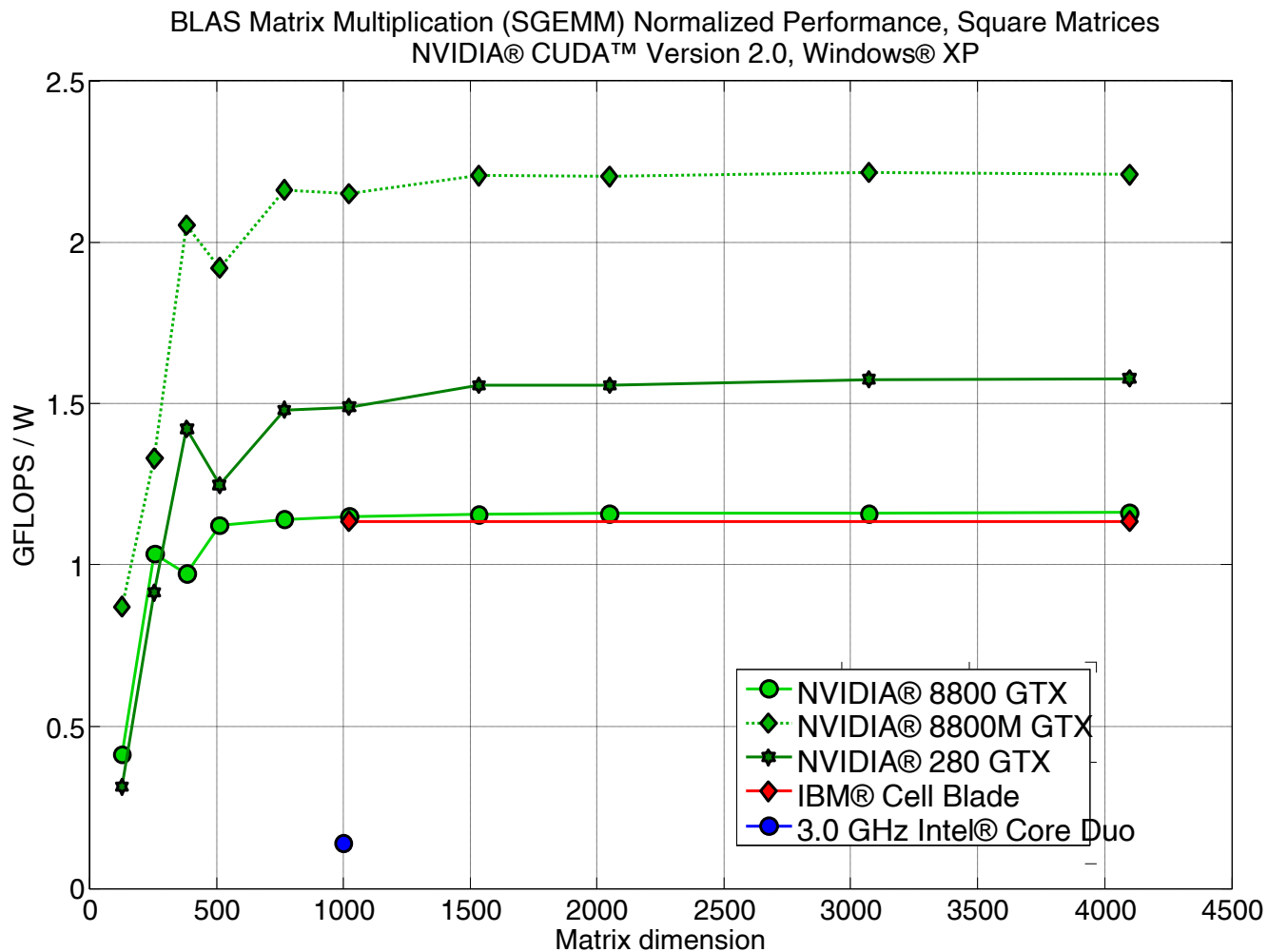
# IRSTAP Benchmarks GFLOP/W

# Matrix Multiplication GFLOP/s

BLAS Matrix Multiplication (SGEMM) Performance, Square Matrices
NVIDIA ® CUDA Version 2.0, Windows XP ®



**GPU Multiplication is SGEMM is from the NVIDIA® library**
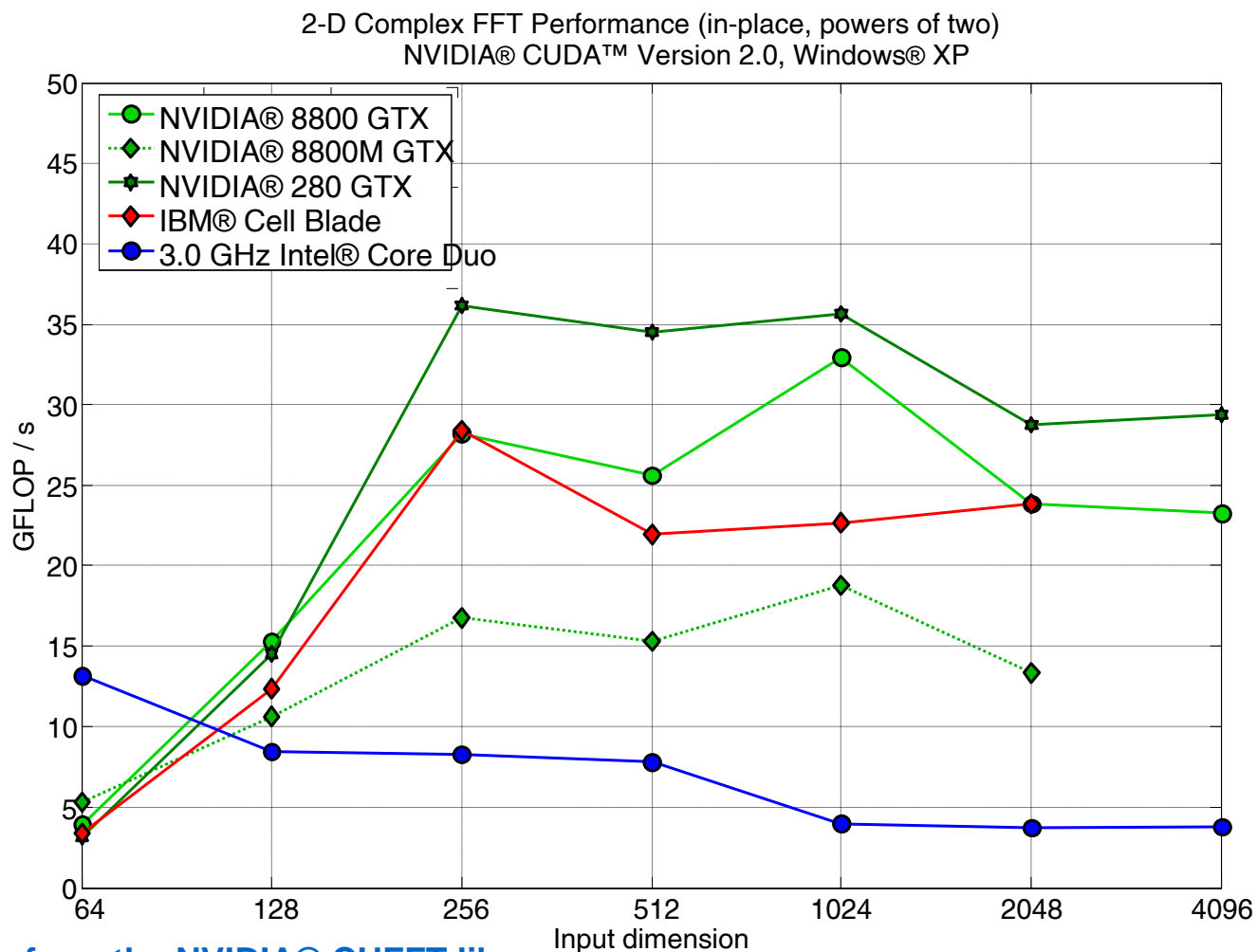
NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries. Intel is a registered trademark of Intel Corporation in the United States and/or other countries.

# Matrix Multiplication GFLOP/W



BLAS Matrix Multiplication (SGEMM) Normalized Performance, Square Matrices
NVIDIA® CUDA™ Version 2.0, Windows® XP

Legend:
- NVIDIA® 8800 GTX
- NVIDIA® 8800M GTX
- NVIDIA® 280 GTX
- IBM® Cell Blade
- 3.0 GHz Intel® Core Duo

Axis labels: GFLOPS / W (vertical), Matrix dimension (horizontal)

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries. Intel is a registered trademark of Intel Corporation in the United States and/or other countries.
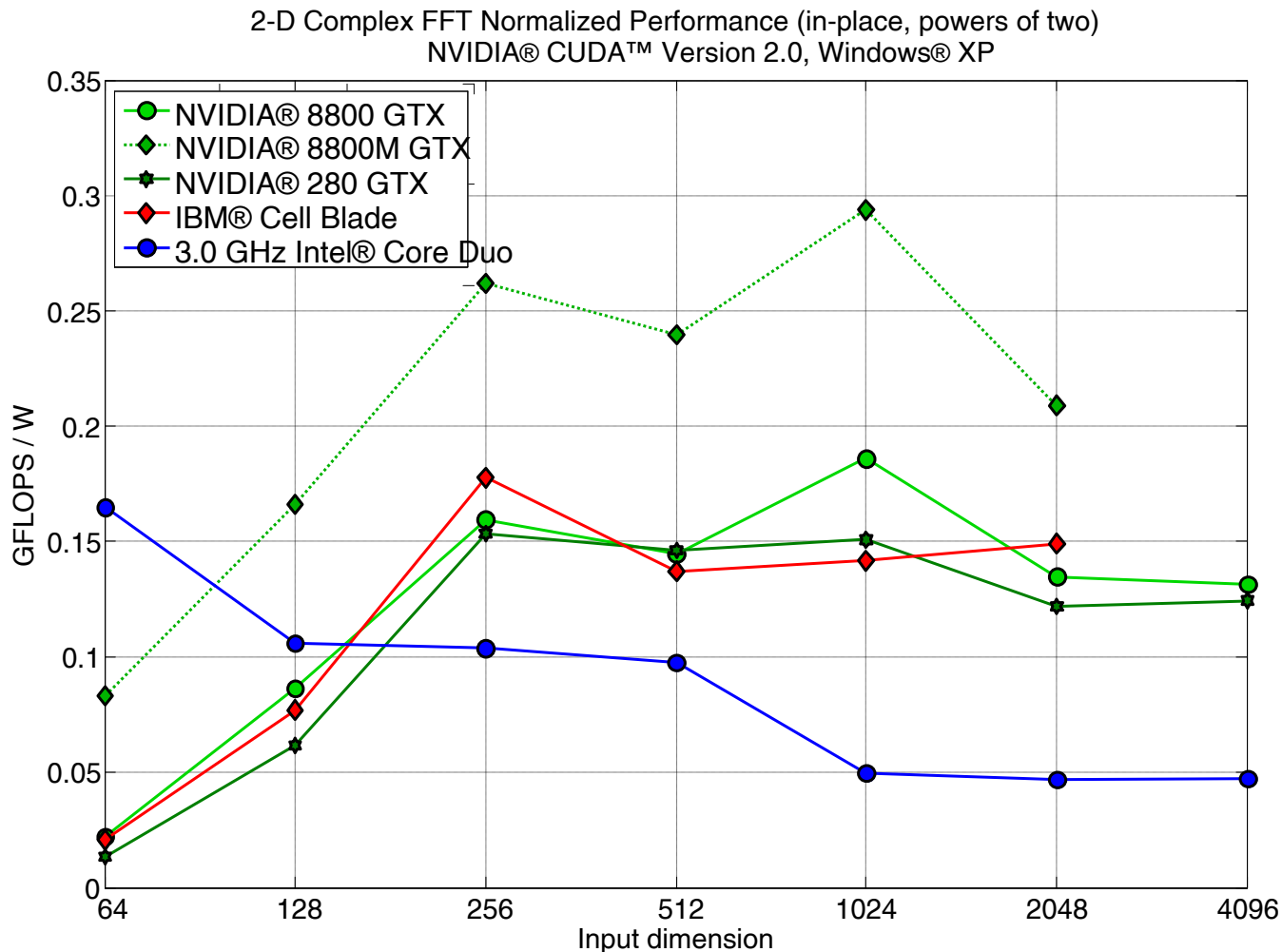
# 2D Complex FFT GFLOP/s



2-D Complex FFT Performance (in-place, powers of two)
NVIDIA® CUDA™ Version 2.0, Windows® XP

Legend:
- NVIDIA® 8800 GTX
- NVIDIA® 8800M GTX
- NVIDIA® 280 GTX
- IBM® Cell Blade
- 3.0 GHz Intel® Core Duo
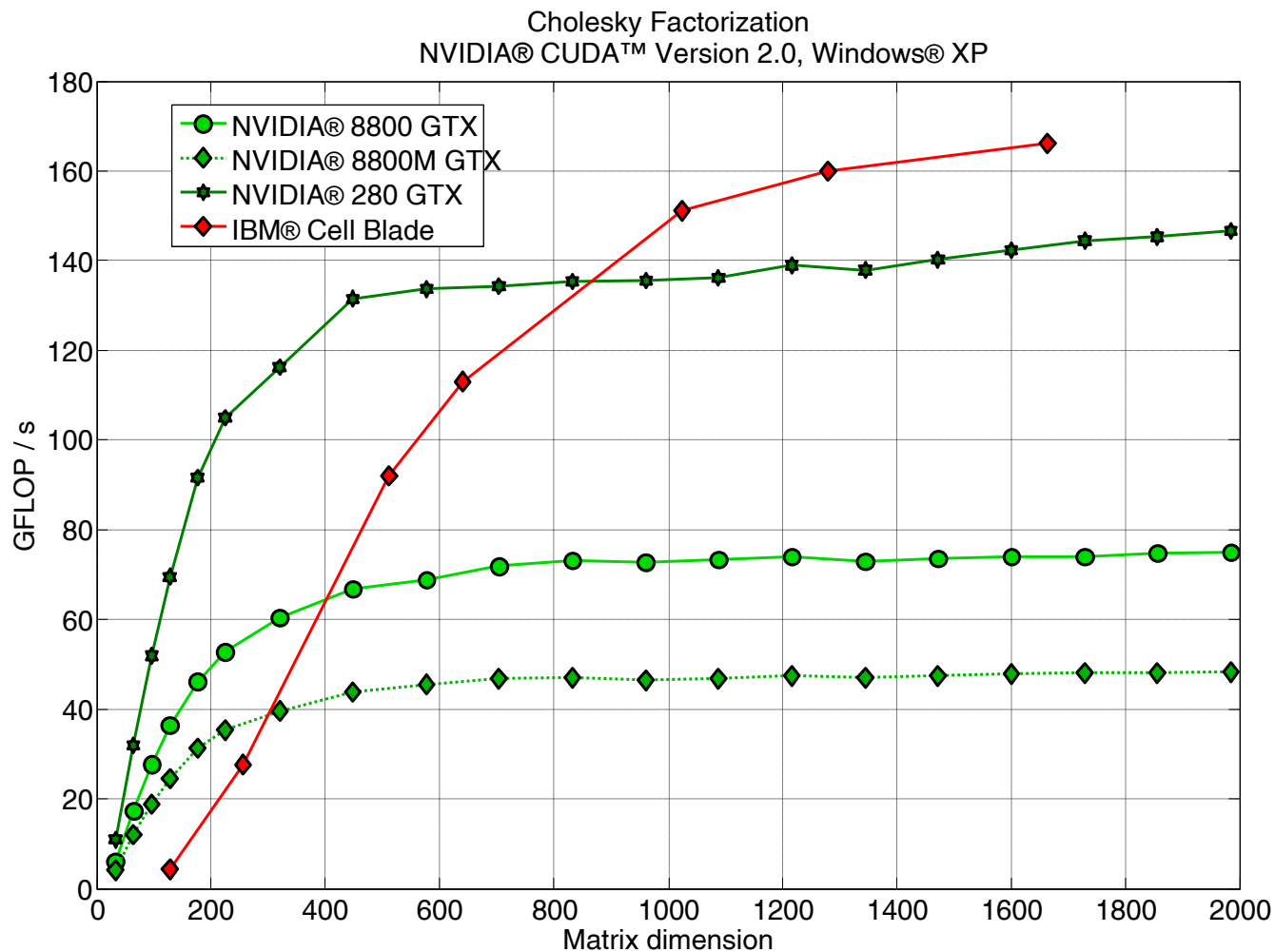
**GPU FFT is from the NVIDIA® CUFFT library**

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries. Intel is a registered trademark of Intel Corporation in the United States and/or other countries.

SAIC
From Science to Solutions

# 2D Complex FFT GFLOP/W



2-D Complex FFT Normalized Performance (in-place, powers of two)
NVIDIA® CUDA™ Version 2.0, Windows® XP

Legend:
- NVIDIA® 8800 GTX
- NVIDIA® 8800M GTX
- NVIDIA® 280 GTX
- IBM® Cell Blade
- 3.0 GHz Intel® Core Duo
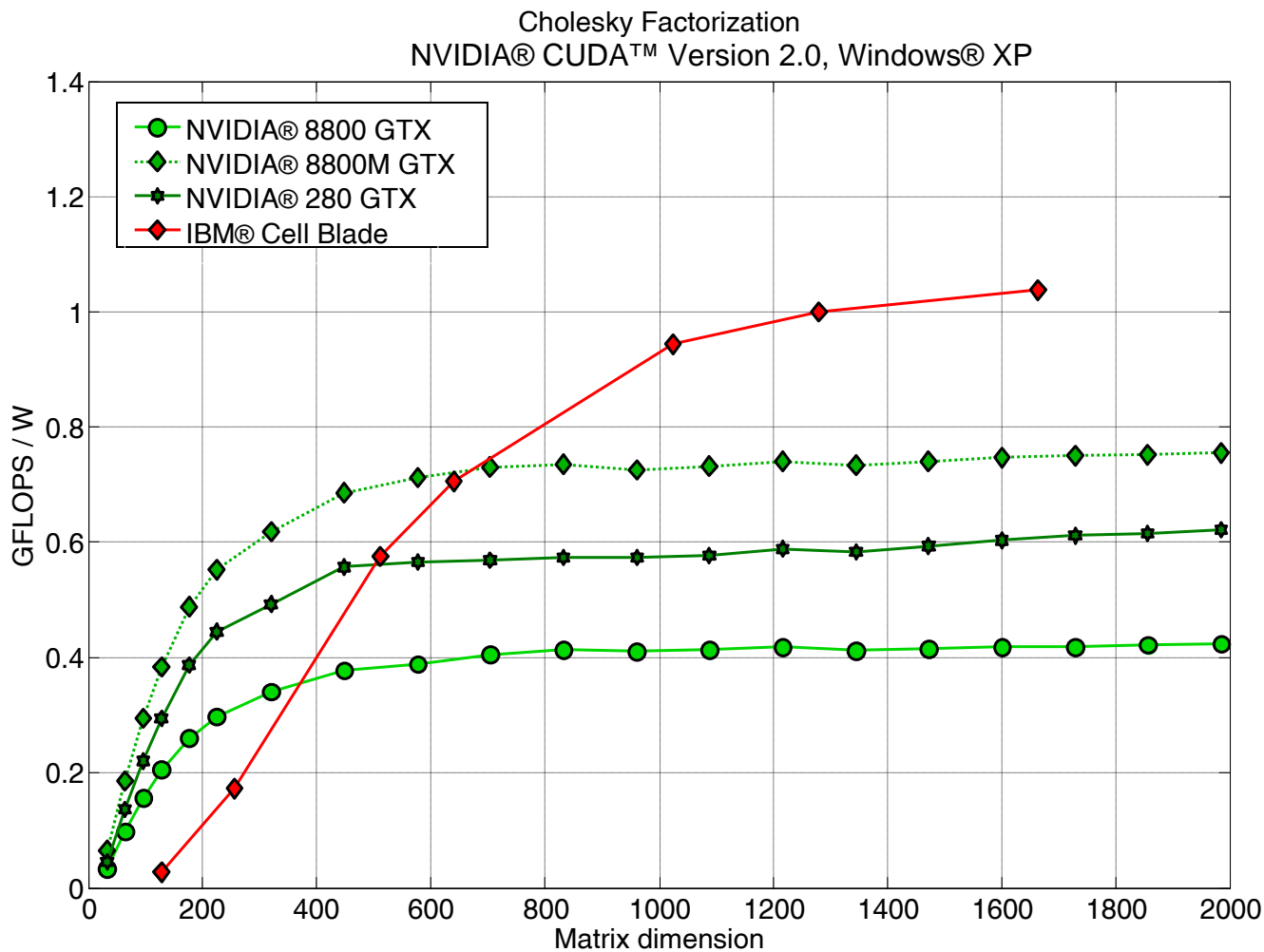
Y-axis: GFLOPS / W
X-axis: Input dimension

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries. Intel is a registered trademark of Intel Corporation in the United States and/or other countries.

# Cholesky GFLOP/s

Cholesky Factorization
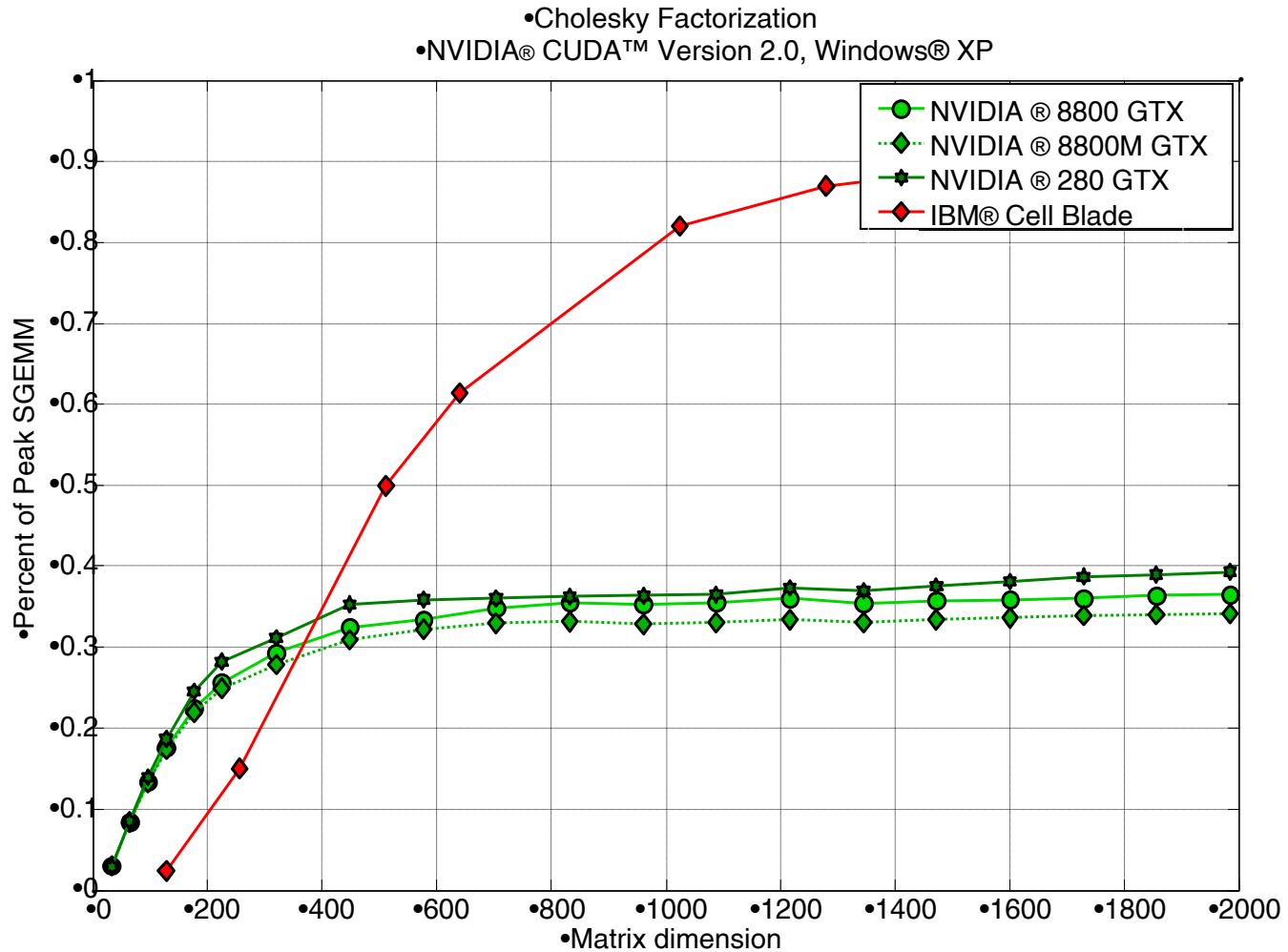NVIDIA® CUDA™ Version 2.0, Windows® XP



NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries.

# Cholesky GFLOP/W



Cholesky Factorization
NVIDIA® CUDA™ Version 2.0, Windows® XP

Legend:
- NVIDIA® 8800 GTX
- NVIDIA® 8800M GTX
- NVIDIA® 280 GTX
- IBM® Cell Blade

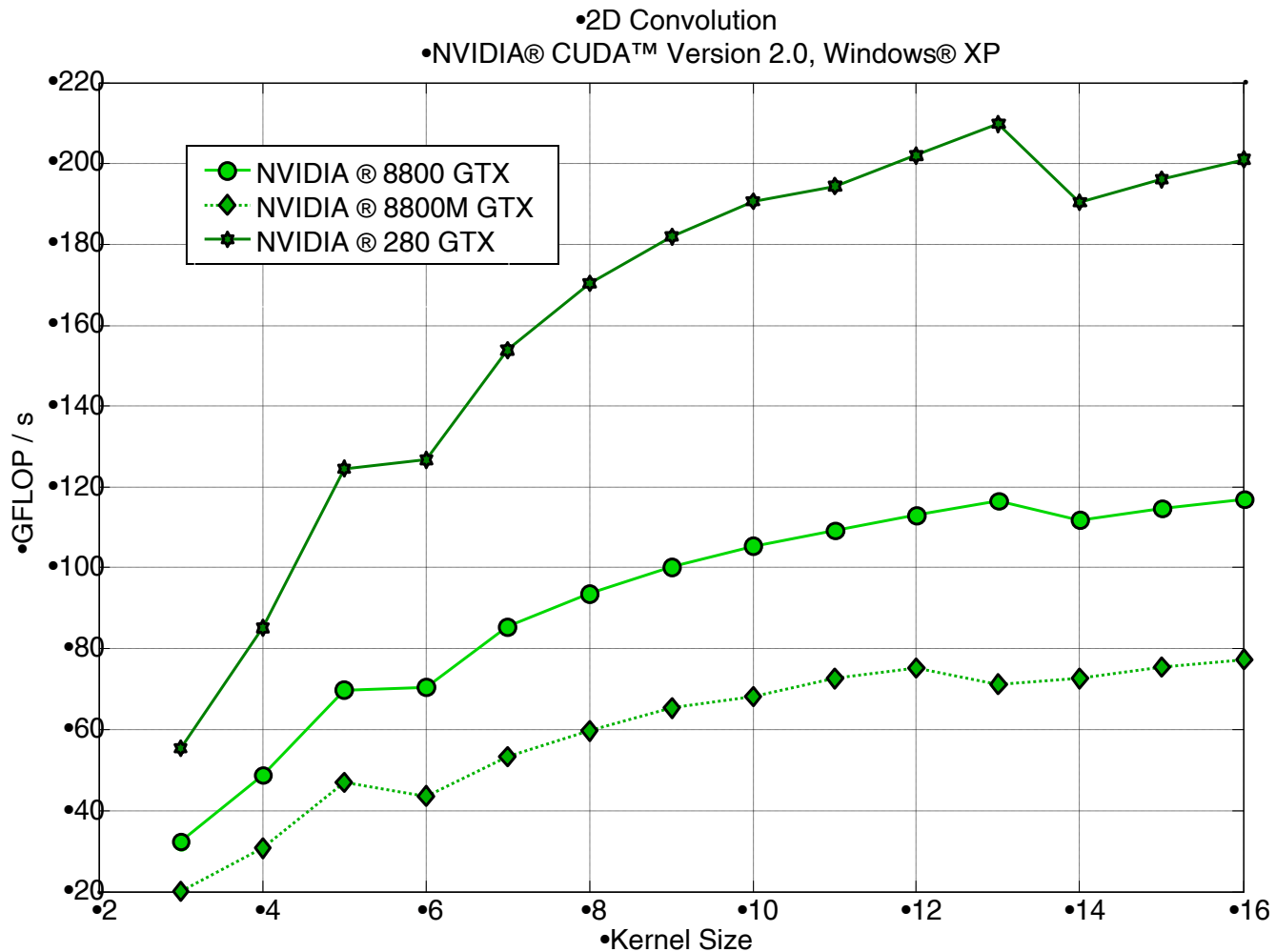Y-axis: GFLOPS / W
X-axis: Matrix dimension

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries.

SAIC
From Science to Solutions

Cholesky Factorization
NVIDIA® CUDA™ Version 2.0, Windows® XP

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries. IBM is a registered trademark of International Business Machines Corporation in the United States and/or other countries.

# 2D Convolution GFLOP/s



•2D Convolution
•NVIDIA® CUDA™ Version 2.0, Windows® XP

Legend:
- NVIDIA ® 8800 GTX
- NVIDIA ® 8800M GTX
- NVIDIA ® 280 GTX
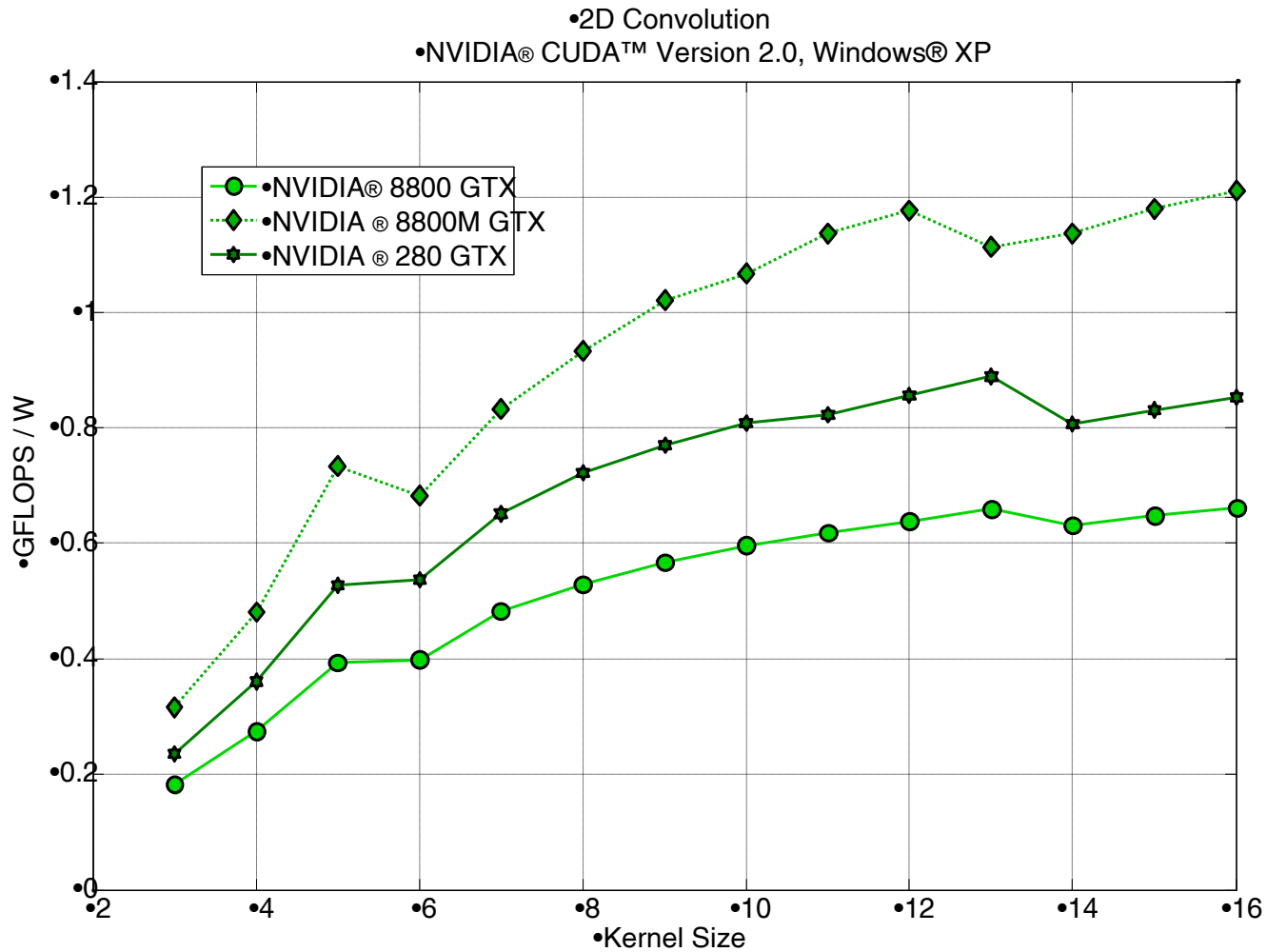
Y-axis: •GFLOP / s
X-axis: •Kernel Size

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries.

SAIC
From Science to Solutions

# 2D Convolution GFLOP/W



•2D Convolution
•NVIDIA® CUDA™ Version 2.0, Windows® XP

Legend:
- •NVIDIA® 8800 GTX
- •NVIDIA ® 8800M GTX
- •NVIDIA ® 280 GTX

Y-axis: •GFLOPS / W (•0, •0.2, •0.4, •0.6, •0.8, •1, •1.2, •1.4)
X-axis: •Kernel Size (•2, •4, •6, •8, •10, •12, •14, •16)

NVIDIA and CUDA are trademarks or registered trademarks of NVIDIA Corporation in the United States and/or other countries. Windows is a registered trademark of Microsoft Corporation in the United States and/or other countries.

18

# Power Estimation

- **Using thermal design power (TDP) for each device**

    – NVDIA® 8800 GTX 128SP 1.35Ghz – TDP: 177W

    – NVDIA® 8800M GTX 96SP 1.25Ghz – TDP: 65W

    – NVDIA® 280 GTX 240SP 1.30Ghz – TDP: 236W

- **Cell – estimating TDP: 160W with 90nm process**

NVIDIA is a registered trademark of NVIDIA Corporation in the United States and/or other countries.

*SAIC*
*From Science to Solutions*

# Questions ?

# Just for comparison…

Assume Cell Processor power consumption is approximately 160 W

2 Cell Processors per blade

---

7000 ft at 32°C

**Power**

Cell Accelerator Board 2 with 4-GB DDR2    162W

Power is provided through the use of a single cable connector in addition to the 75W power provided through the PCI Express edge connector.

---

- Power consumption:
  - QS20: 315 watts maximum
  - QS20: 330 watts with 1 IB (#2945)
  - QS20: 345 watts with 2 IB (#2945)

*SAIC*
From Science to Solutions