# Accelerated Computing from Mobile Devices to Supercomputers
*Dale Southard, NVIDIA*

Power of 600 Petaflop
CPU-only Supercomputer = Power for the city
of San Francisco

## HPC's Biggest Challenge: Power

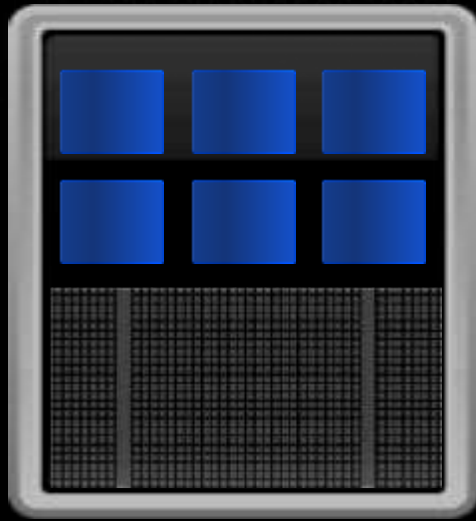# GPUs Power World's 10 Greenest Supercomputers

| Green500 Rank | MFLOPS/W | Site |
|---|---|---|
| 1 | 4,503.17 | GSIC Center, Tokyo Tech |
| 2 | 3,631.86 | Cambridge University |
| 3 | 3,517.84 | University of Tsukuba |
| 4 | 3,185.91 | Swiss National Supercomputing (CSCS) |
| 5 | 3,130.95 | ROMEO HPC Center |
| 6 | 3,068.71 | GSIC Center, Tokyo Tech |
| 7 | 2,702.16 | University of Arizona |
| 8 | 2,629.10 | Max-Planck |
| 9 | 2,629.10 | (Financial Institution) |
| 10 | 2,358.69 | CSIRO |
| 37 | 1959.90 | Intel Endeavor (top Xeon Phi cluster) |
| 49 | 1247.57 | Météo France (top CPU cluster) |

THE GREEN 500™

# Accelerated Computing
## *10x Performance & 5x Energy Efficiency for HPC*

**GPU Accelerator**
Optimized for
Parallel Tasks

**CPU**
Optimized for
Serial Tasks

+

# How GPU Acceleration Works

**Application Code**

Compute-Intensive Functions

5% of Code

Rest of Sequential
CPU Code

**GPU**

**CPU**

+

# Accelerated Computing Growing Fast

## 2x Growth in One Year

**Percent of HPC Systems With Accelerators**

- 50%
- 40%
- 30%
- 20%
- 10%
- 0%

44%

22%  24%

2010  2011  2012

## Hundreds of GPU Accelerated Apps

- 300
- 250
- 200
- 150
- 100
- 50
- 0

242

182

113

2011  2012  2013

## NVIDIA GPU is Accelerator of Choice

INTEL PHI
4%

OTHERS
11%

NVIDIA GPUs
85%

NVIDIA

POPULAR GPU-ACCELERATED APPLICATIONS

## Research: Higher Education and Supercomputing

### COMPUTATIONAL CHEMISTRY AND BIOLOGY

272 GPU-Accelerated Applications
www.nvidia.com/appscatalog

# Artificial Neural Network at a Fraction of the Cost with GPUs

GOOGLE BRAIN

1,000 CPU Servers
2,000 CPUs • 16,000 cores

**600 kWatts**
**$5,000,000**

STANFORD AI LAB

3 GPU-Accelerated Servers
12 GPUs • 18,432 cores

**4 kWatts**
**$33,000**

*" Now You Can Build Google's $1M Artificial Brain on the Cheap "*

-Wired

# GPUs Accelerate Machine Learning & Data Analytics


Auto Tagging in Creative Cloud


Speech/Image Recognition


Hadoop-based Clustering


Recommendation Engine


Database Queries


Search Ranking

# IBM Partners with NVIDIA to Build Next-Generation Supercomputers

Tesla
**GPU**

**+**

POWER 8
**CPU**

GPU-Accelerated POWER-Based Systems Available in 2014

# JETSON TK1
## THE WORLD'S 1st EMBEDDED SUPERCOMPUTER

Development Platform for Embedded
Computer Vision, Robotics, Medical

192 Cores · 326 GFLOPS FP32

CUDA Enabled

Available Now

# GPU Accelerated Libraries
## "Drop-in" Acceleration for your Applications

**Linear Algebra**
FFT, BLAS,
SPARSE, Matrix


NVIDIA cuFFT, cuBLAS, cuSPARSE


CULA|tools


MAGMA


CUSP

**Numerical & Math**
RAND, Statistics


IMSL Fortran Numerical Library


NVIDIA Math Lib


ArrayFire


NVIDIA cuRAND

**Data Struct. & AI**
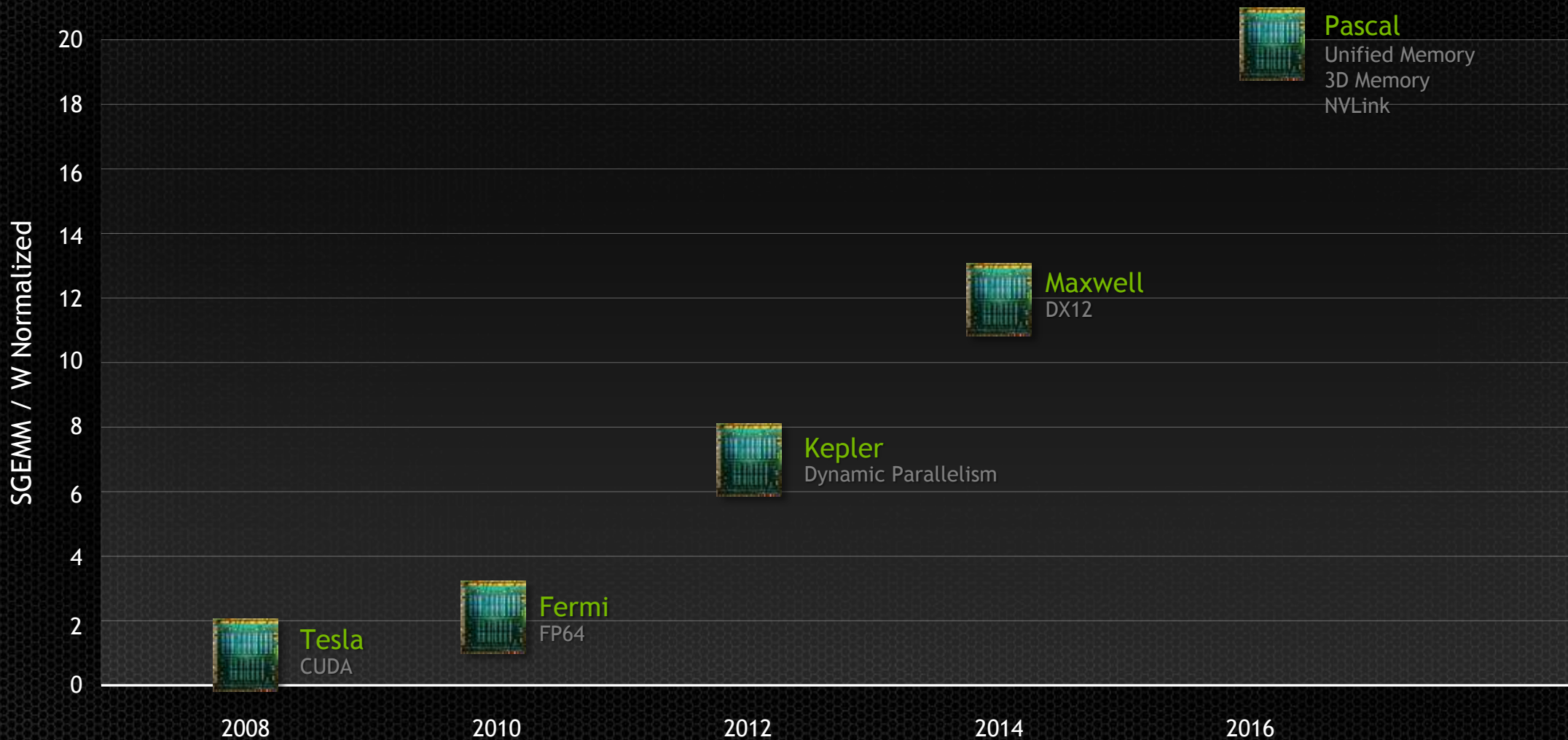Sort, Scan, Zero Sum


Thrust


GPU AI – Board Games


GPU AI – Path Finding

**Visual Processing**
Image & Video


NVIDIA NPP


NVIDIA Video Encode


Sundog Software

# OpenACC: Open, Simple, Portable

```
main() {

  ...
  <serial code>

  ...
  #pragma acc kernels

  {
  <compute intensive code>

  }

  ...
}
```

**Compiler Hint**

- Open Standard
- Easy, Compiler-Driven Approach
- Portable on GPUs and Xeon Phi

**CAM-SE Climate**
6x Faster on GPU
Top Kernel: 50% of Runtime

# Linux GCC Compiler to Support GPU Accelerators

## Open Source
GCC Efforts by Samsung & Mentor Graphics

## Pervasive Impact
Free to all Linux users

## Mainstream
Most Widely Used HPC Compiler

> *" Incorporating OpenACC into GCC is an excellent example of open source and open standards working together to make accelerated computing broadly accessible to all Linux developers. "*
>
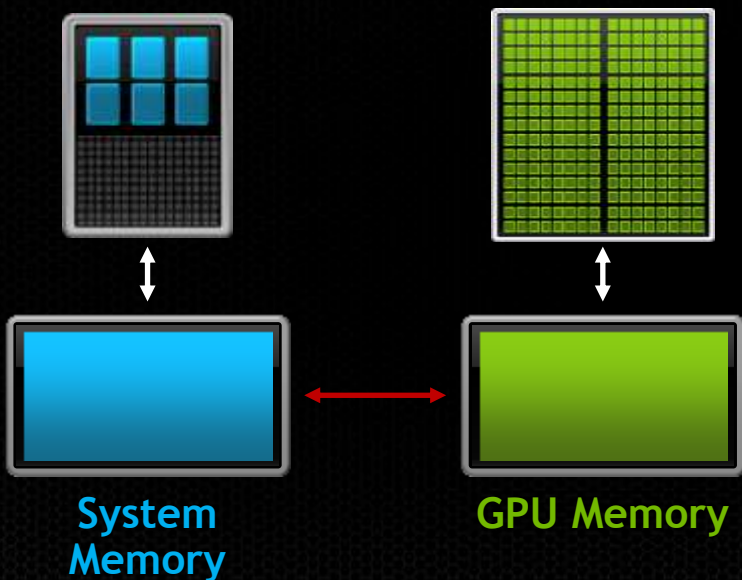> Oscar Hernandez
> Oak Ridge National Laboratories

OpenACC
Directives for Accelerators

# Unified Memory
## Dramatically Lower Developer Effort

**Developer View Today**

**Developer View With Unified Memory**

**System Memory**

**GPU Memory**

**Unified Memory**

# CUDA: World's Most Pervasive Parallel Programming Model

**14,000** — Institutions with CUDA Developers

**2,000,000** — CUDA Downloads

**487,000,000** — CUDA GPUs Shipped

**700+ University Courses In 62 Countries**

**Dale Southard, dsouthard@nvidia.com**
Principal System Architect, Office of the CTO