*GPUs and the Future of Accelerated Computing*
Emerging Technology Conference 2014
University of Manchester

**John Ashley**
*Senior Solutions Architect*
*jashley@nvidia.com*

# Disclaimers

- These are my views, not NVIDIA's (although I do hope you all will agree by the time I'm done!) , and where I've quoted others, any misrepresentation of their views is mine and mine alone.

- There are "forward looking statements" in here. You are warned that my crystal ball is no better than anyone else's, even if I'm speaking about things like NVIDIA roadmaps, etc. It's a lawyer thing.

- Trademarks of other firms are their own, etc.

# Agenda

- **Who is NVIDIA?**

- **Why should I care about accelerated computing?**

- **What sort of difference can CUDA make?**

- **Where is NVIDIA going?**

# Who is NVIDIA?

NVIDIA products span the power-performance spectrum

# 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) |

# GPUs are becoming…

100M
CUDA –Capable
GPUs

150K
CUDA Downloads

77
Supercomputing
Teraflops

60
University
Courses

4,000
Academic Papers
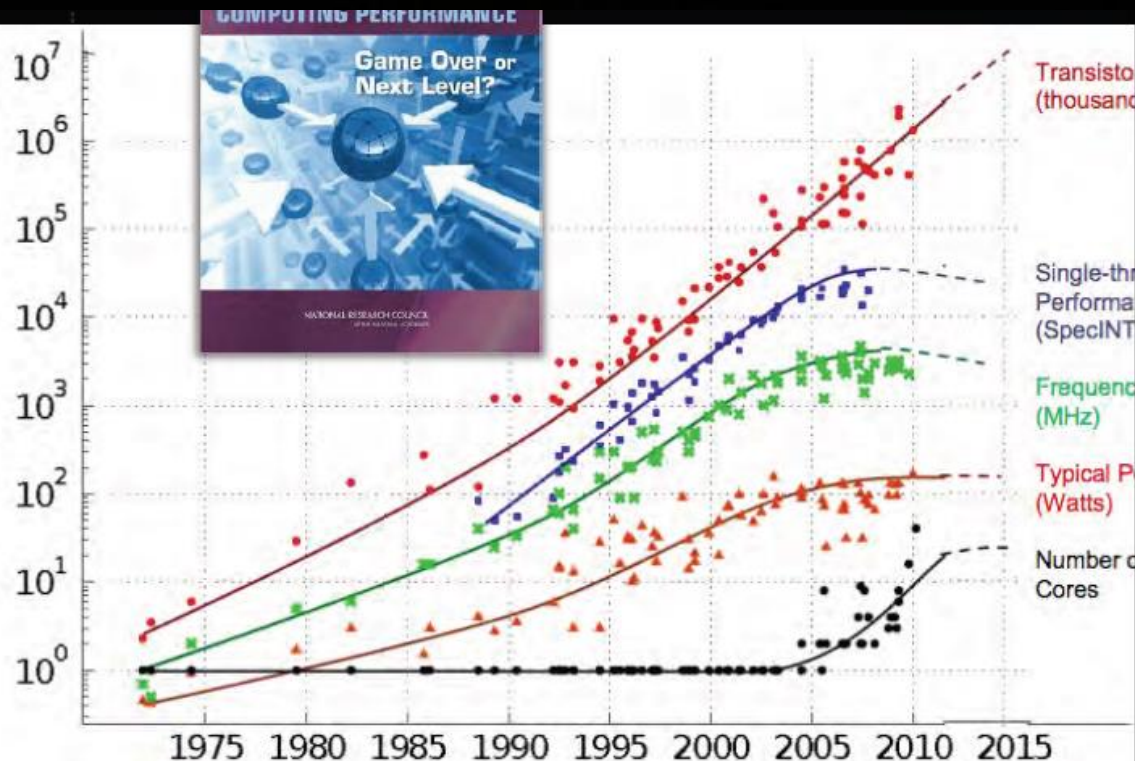
2008    2014

# GPUs are becoming…pervasive

**100M** CUDA –Capable GPUs

**150K** CUDA Downloads

**77** Supercomputing Teraflops

**60** University Courses

**4,000** Academic Papers

**430M** CUDA-Capable GPUs

**2.2M** CUDA Downloads

**41,700** Supercomputing Teraflops

**738** University Courses

**50,000** Academic Papers

2008    2014

**Why should I care about accelerated computing?**

# Moore's Law isn't what it used to be.



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Ba...
Dotted line extrapolations by C. Moore

## Moore's law is alive and well, but...

Instruction-level parallelism (ILP) was mined out in 2001

Voltage scaling (Dennard scaling) ended in 2005

Most power is spent on communication

## What does this mean to you?

# The future is here today …

- Performance gains come from parallelism

- Systems are power limited
  (efficiency IS performance)

- Systems are communication
  limited
  (locality IS performance)

# High Performance Computing is Going Hybrid



PCIe

**CPU***
Fast single threads
(serial work)

Sandy Bridge
32nm
690 pJ/flop

**GPU**
Extreme power-efficiency
(throughput work)

Kepler
28nm
134 pJ/flop

- Do most work by many cores optimized for **extreme energy efficiency**

- Still need a few cores optimized for **fast serial work**

- Amdahl's law continues to apply to all architectures

*x86, ARM, Power

12

# What sort of difference can CUDA make?

# CUDA Parallel Computing Platform

www.nvidia.com/getcuda

**Programming Approaches**

| Libraries | Directives | Programming Languages |
|---|---|---|
| "Drop-in" Acceleration | Easily Accelerate Apps | Maximum Flexibility |

**Development Environment**

Nsight IDE
Linux, Mac and Windows
GPU Debugging and Profiling

CUDA-GDB debugger
NVIDIA Visual Profiler

**Open Compiler Tool Chain**

LLVM COMPILER INFRASTRUCTURE

Enables compiling new languages to CUDA platform, and CUDA languages to other architectures

**Hardware Capabilities**

SMX          Dynamic Parallelism          HyperQ          GPUDirect

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

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

## 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**

# VISIONWORKS
# COMPUTER VISION ON CUDA

Driver Assistance

Computational Photography

Augmented Reality

Robotics

**Your Code**

**Sample Pipelines**

Object Detection / Tracking

Structure from Motion

. . .

**VisionWorks Primitives**

Classifier

Corner Detection

. . .

**CUDA**

**Jetson TK1**

# Computer Vision on CUDA



Feature Detection / Tracking

~30 GFLOPS @ 30 Hz



3D Scene Interpretation

~280  GFLOPS @ 30 Hz

# Audi Self Driving Car Before & After CUDA

# Solid Growth of GPU Accelerated Apps

**# of GPU-Accelerated Apps**



Bar chart values:
- 2011: 113
- 2012: 182
- 2013: 272

## Top HPC Applications

| Category | | |
|---|---|---|
| Molecular Dynamics | AMBER<br>CHARMM<br>DESMOND | GROMACS<br>LAMMPS<br>NAMD |
| Quantum Chemistry | Abinit<br>Gaussian | GAMESS<br>NWChem |
| Material Science | CP2K<br>QMCPACK | Quantum Espresso<br>VASP |
| Weather & Climate | COSMO<br>GEOS-5<br>HOMME | CAM-SE<br>NEMO<br>NIM<br>WRF |
| Lattice QCD | Chroma | MILC |
| Plasma Physics | GTC | GTS |
| Structural Mechanics | ANSYS Mechanical<br>LS-DYNA Implicit<br>MSC Nastran | OptiStruct<br>Abaqus/Standard |
| Fluid Dynamics | ANSYS Fluent | Culises<br>(OpenFOAM) |

Accelerated, In Development

19

POPULAR GPU-ACCELERATED APPLICATIONS

## Research: Higher Education and Supercomputing

### COMPUTATIONAL CHEMISTRY AND BIOLOGY

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

# Top Applications Now with Built-in GPU Support

**Digital Content Creation**

Adobe CS

Apple Final Cut

Non-GPU Apps

Autodesk 3dsMax

Sony Vegas Pro

Other GPU Apps

Avid Media Composer

**Application Market Share by Segment**

**Molecular Dynamics**

NAMD

AMBER

GROMACS

Non-GPU Apps

CHARMM

DL_POLY

LAMMPS

**Computer-Aided Engineering**

ANSYS

Simulia Abaqus

Non-GPU Apps

MSC Nastran

Altair Radioss

**Quantum Chemistry**

GAMESS

Gaussian

NWChem

Non-GPU Apps

Quantum Espresso

CP2K

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

# Performance on Leading Scientific Applications

# Where is NVIDIA going?

# Fast Paced CUDA GPU Roadmap

Volta

Pascal

Unified Memory
Stacked DRAM
NVLINK

Maxwell

Higher Perf/Watt

Kepler

Dynamic Parallelism

Fermi

FP64

Tesla

CUDA

GFLOPS per Watt

32
16
8
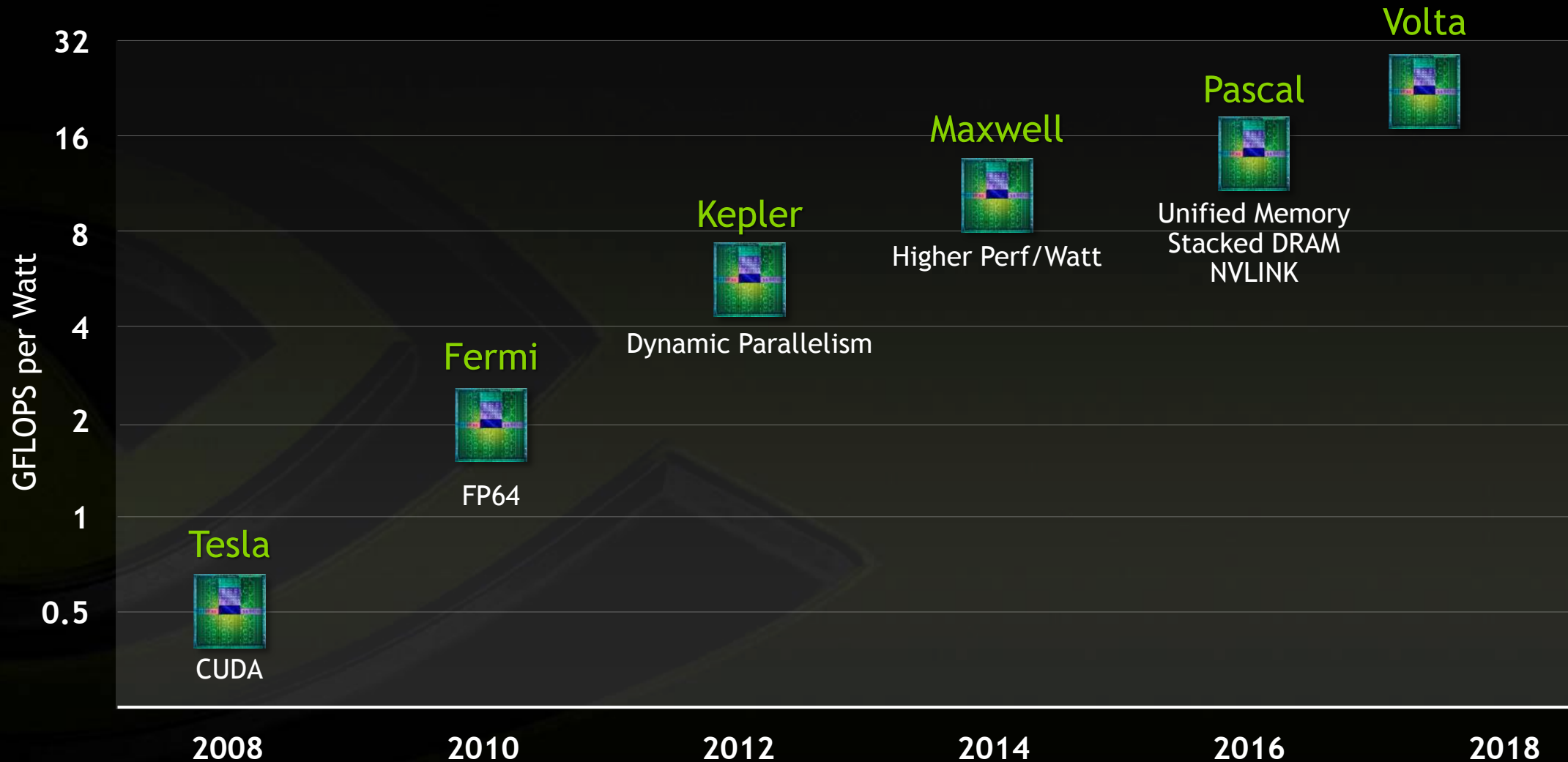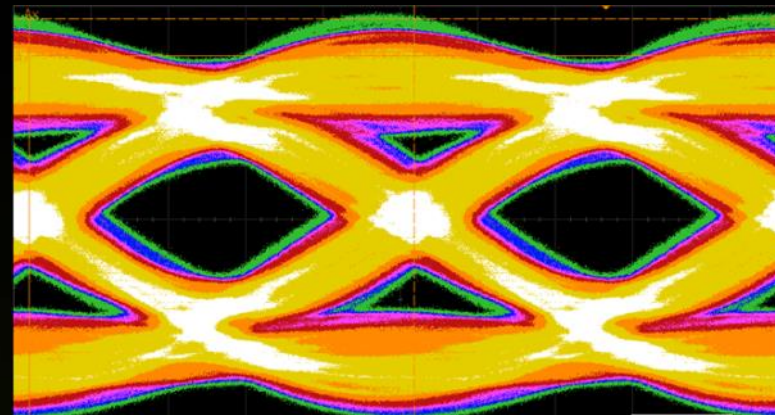4
2
1
0.5

2008   2010   2012   2014   2016   2018

# Introducing NVLINK and Stacked Memory
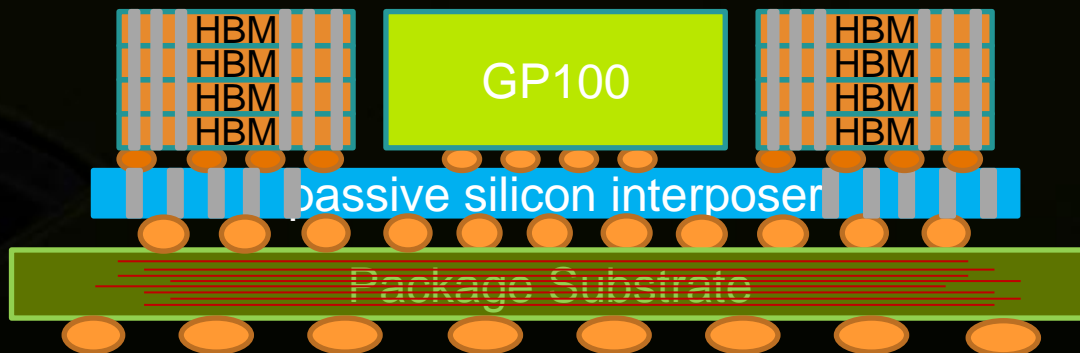
## NVLINK

- **GPU high speed interconnect**
- **5-12x PCIe Gen 3 Bandwidth**
- **Drastically reduced energy/bit**

## Stacked Memory
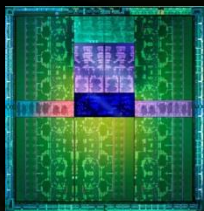
- 2-4x Capacity & Bandwidth
- 3-4x More Energy Efficient per bit
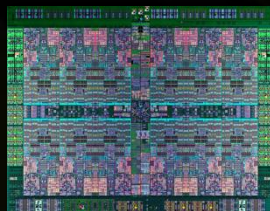- Leaves more power for compute

# IBM Partners with NVIDIA to Build Next-Generation Supercomputers
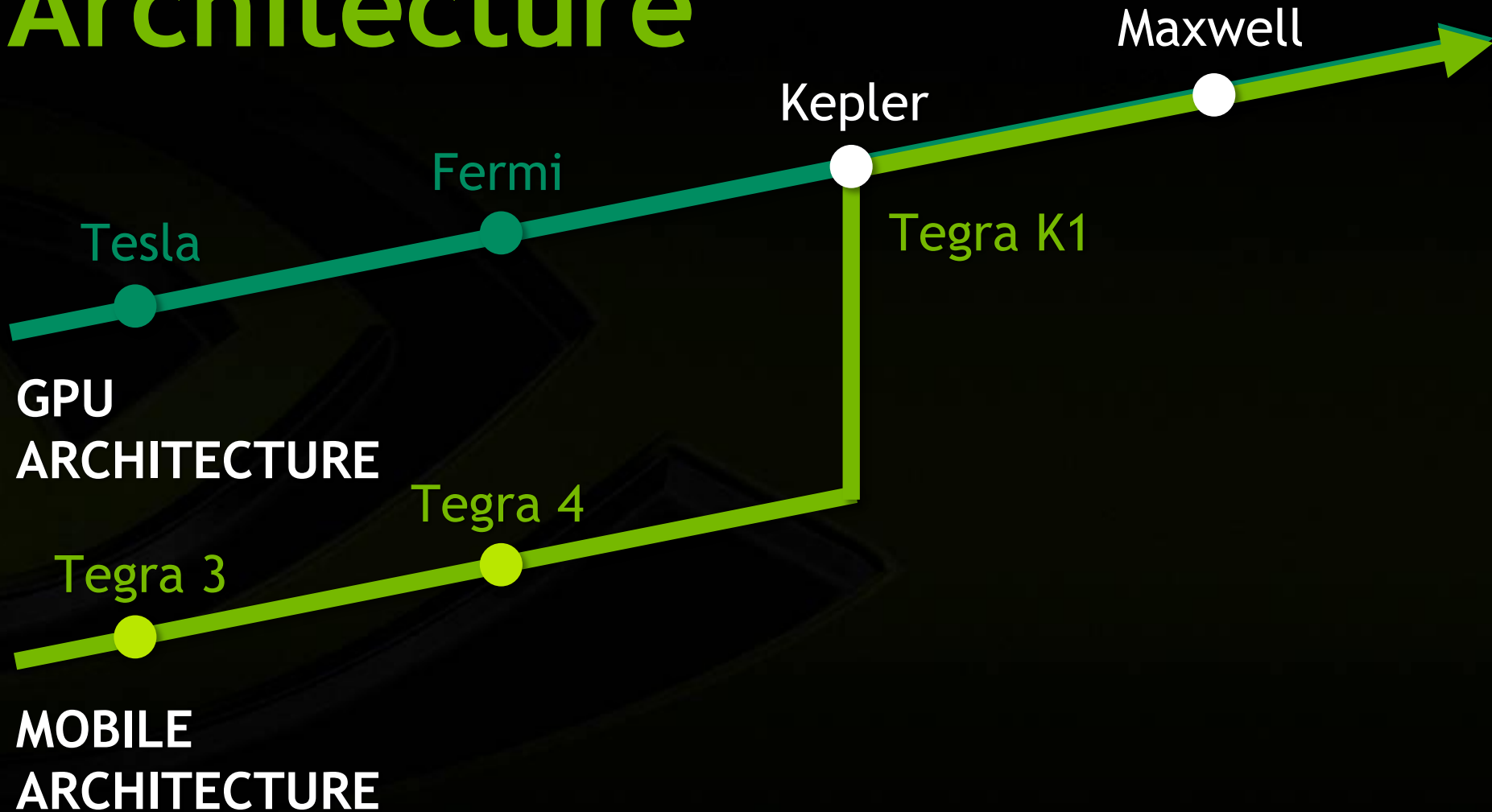


Tesla
GPU

+

POWER 8
CPU

GPU-Accelerated POWER-Based Systems Available in 2014

# Unify GPU and Tegra Architecture



Maxwell

Kepler

Fermi

Tesla

Tegra K1

**GPU ARCHITECTURE**

Tegra 4

Tegra 3

**MOBILE ARCHITECTURE**

27

# Jetson TK1 Development Kit



- **32 Bit ARM+Kepler SMX SOC**

- **CUDA capable**

- **$192 from US retailers, cost and availability will vary elsewhere**

- **https://developer.nvidia.com/jetson-tk1**

- **http://devblogs.nvidia.com/parallelforall/jetson-tk1-mobile-embedded-supercomputer-cuda-everywhere/**

# Summary

- **Who is NVIDIA?**
  A: The world's leading Visual Computing company, from consumer devices through to world class supercomputers

- **Why should I care about accelerated computing?**
  A: Because parallelism and heterogeneous computing is the future of big compute and big data

- **What sort of difference can CUDA make?**
  A: Order of magnitude improvements in performance and efficiency are possible with CUDA

- **Where is NVIDIA going?**
  A: Relentlessly forward to a CUDA enabled parallel & energy efficient computing future

**Questions ?**