# DEHPC '15

Dale Southard

**NVIDIA.**

GAMING

DESIGN

DESKTOP VIRTUALIZATION

HPC & CLOUD SERVICE PROVIDERS

AUTONOMOUS MACHINES

PC

DATA CENTER

MOBILE

# The World Leader in Visual Computing

NVIDIA.

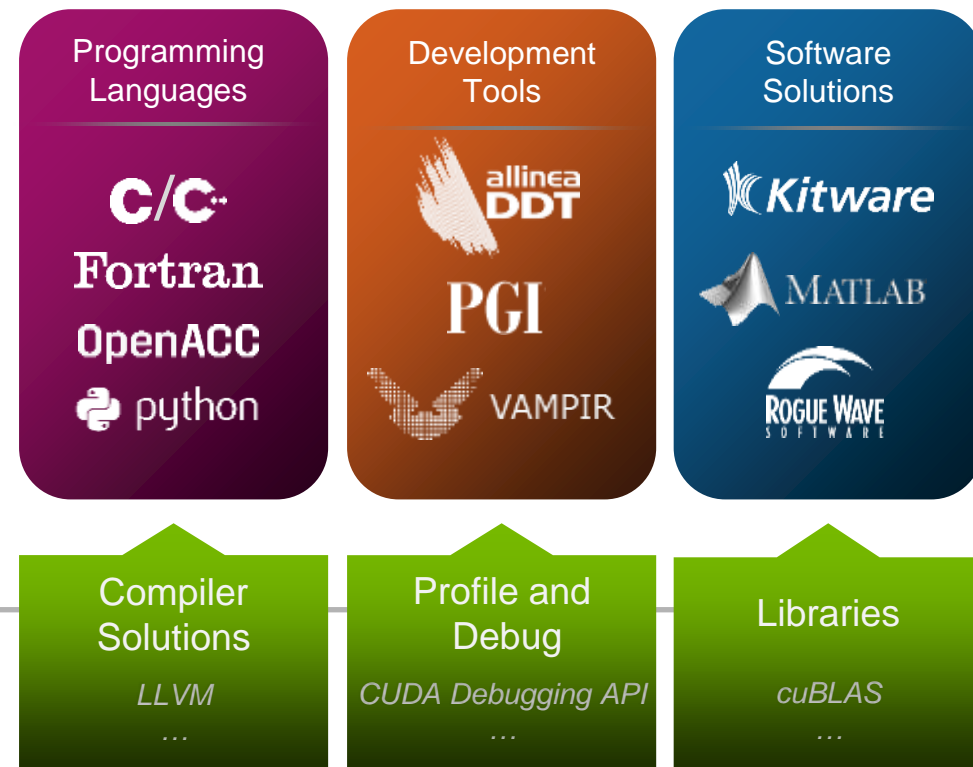# Tesla Accelerated Computing Platform

## Data Center Infrastructure

### System Solutions
CRAY · cisco
DELL · hp
IBM · lenovo
quanta sgi · SUPERMICRO
amazon web services | EC2

### Communication
Mellanox TECHNOLOGIES
MVAPICH
OPEN MPI

### Infrastructure Management
Adaptive COMPUTING
Bright Computing
IBM PLATFORM COMPUTING

#### GPU Accelerators
*GPU Boost*
...

#### Interconnect
*GPU Direct NVLink*
...

#### System Management
*NVML*
...

## Development

### Programming Languages
C/C++
Fortran
OpenACC
python

### Development Tools
allinea DDT
PGI
VAMPIR

### Software Solutions
Kitware
MATLAB
ROGUE WAVE SOFTWARE

#### Compiler Solutions
*LLVM*
...

#### Profile and Debug
*CUDA Debugging API*
...

#### Libraries
*cuBLAS*
...

*" Accelerators Will Be Installed in More than Half of New Systems "*

*"In 2014, NVIDIA enjoyed a dominant market share with 85% of the accelerator market."*

Source: Top 6 predictions for HPC in 2015    Intersect360 RESEARCH

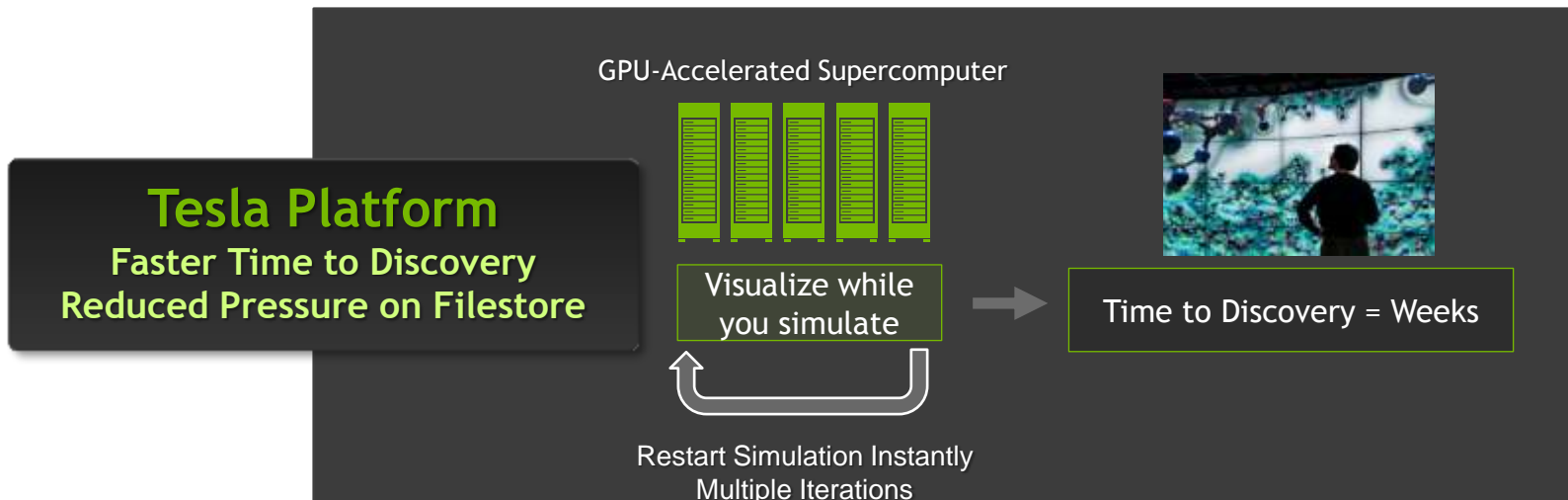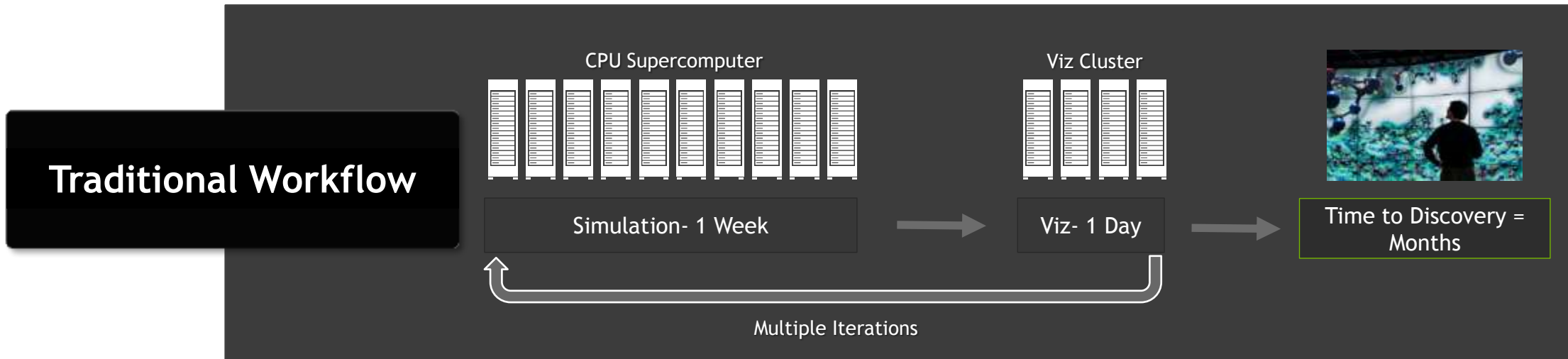3   NVIDIA

# Vision: Mainstream Parallel Programming

Enable more programmers to write portable parallel software in their language of choice
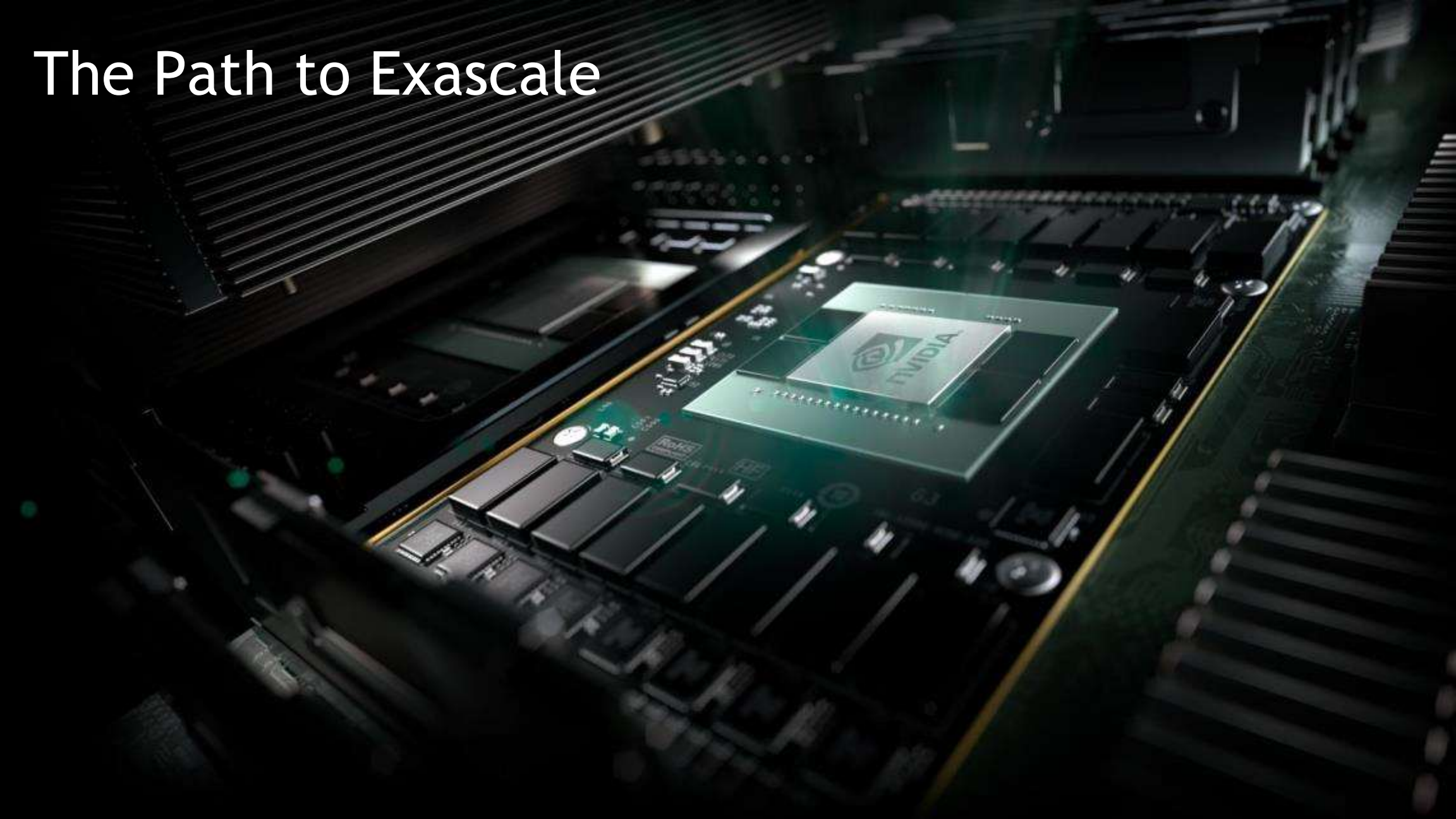
Embrace and evolve standards in key languages

CUDA continues to evolve as the target low-level platform for GPU acceleration

# Vision: In Situ Vis – Faster Science, Lower Cost

**Traditional Workflow**

CPU Supercomputer

Viz Cluster

Simulation- 1 Week → Viz- 1 Day → Time to Discovery = Months

Multiple Iterations

**Tesla Platform**
Faster Time to Discovery
Reduced Pressure on Filestore

GPU-Accelerated Supercomputer

Visualize while you simulate → Time to Discovery = Weeks

Restart Simulation Instantly
Multiple Iterations

⊙ nVIDIA.

# The Path to Exascale
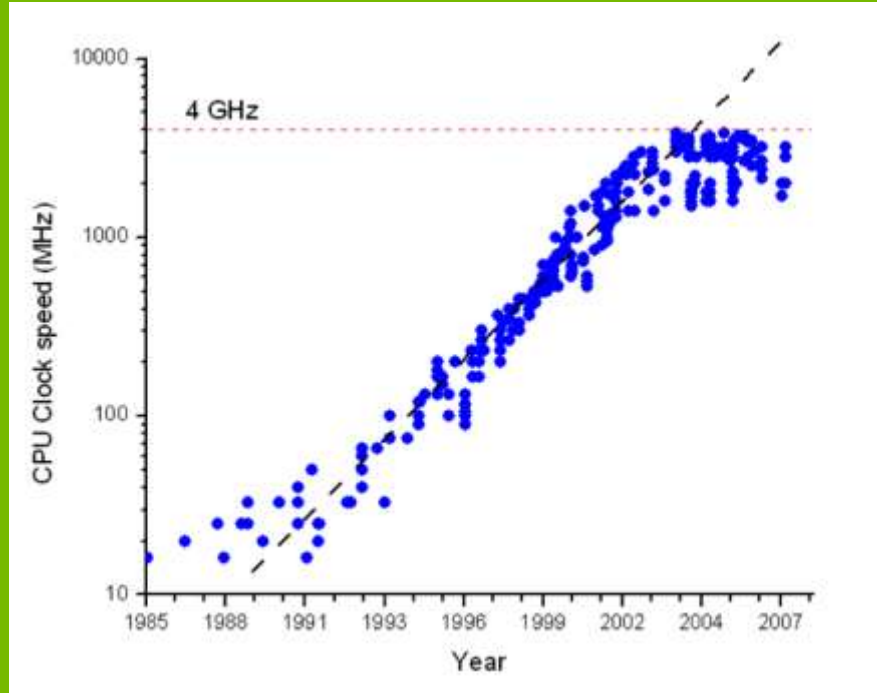
Power for CPU-only **Exaflop** Supercomputer = Power for the Bay Area, CA *(San Francisco + San Jose)*

# HPC's Biggest Challenge

# Hitting a Frequency Wall?



G Bell, *History of Supercomputers*, LLNL, April 2013

# The End of Voltage Scaling

## The Good Old Days

Leakage was not important, and voltage scaled with feature size

$L' = L/2$
$V' = V/2$
$E' = CV^2 = E/8$
$f' = 2f$
$D' = 1/L^2 = 4D$
$P' = P$

Halve L and get 4x the transistors and 8x the capability for the same power

## The New Reality

Leakage has limited threshold voltage, largely ending voltage scaling

$L' = L/2$
$V' = \sim V$
$E' = CV^2 = E/2$
$f' = 2f$
$D' = 1/L2 = 4D$
$P' = 4P$
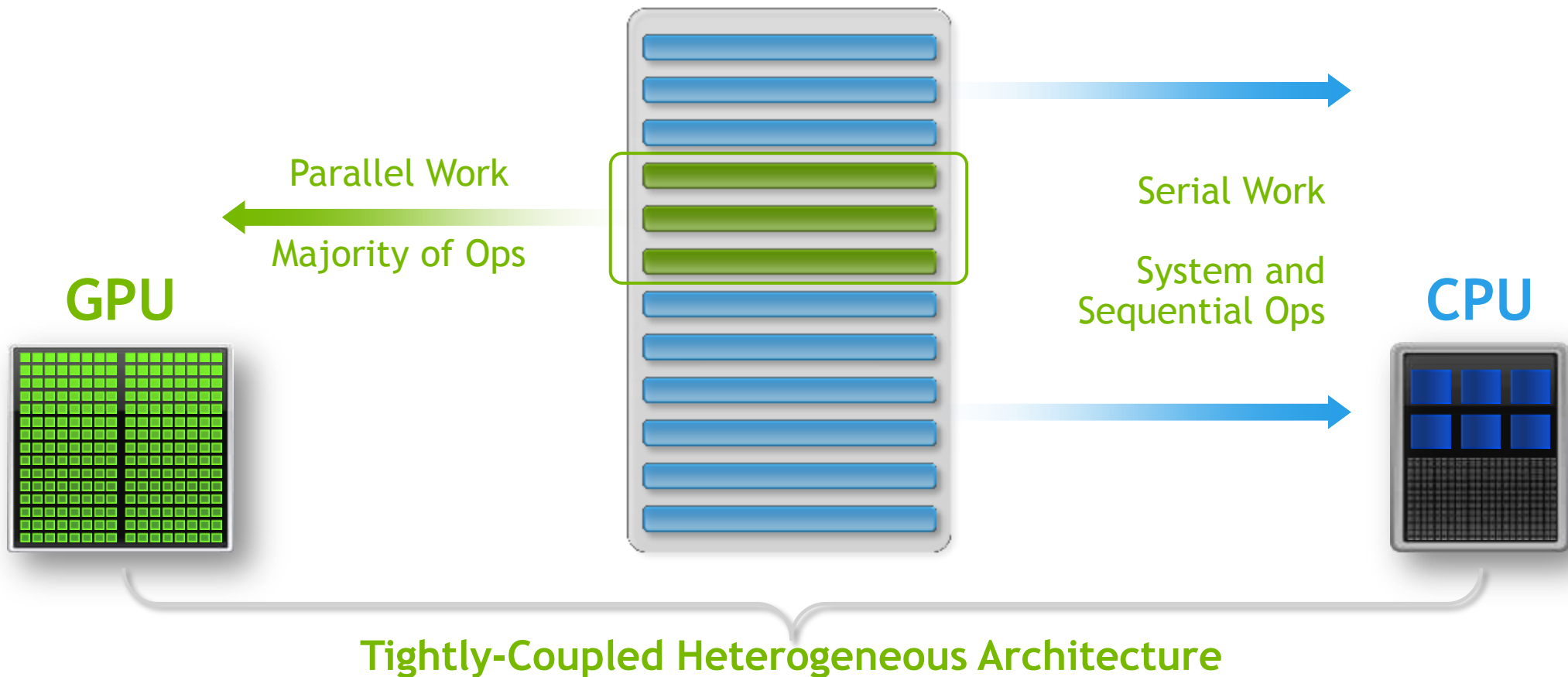
Halve L and get 4x the transistors and 8x the capability for 4x the power, or 2x the capability for the same power in ¼ the area.

"If you want to plow a field, which would you rather use, 4 strong oxen or 1024 chickens?"

*- Seymour Cray, 1989*

**Hint: We want <u>both</u>.**

NVIDIA.

# Optimizing Serial/Parallel Execution

**GPU**

Parallel Work

Majority of Ops

**Tightly-Coupled Heterogeneous Architecture**

Serial Work

System and
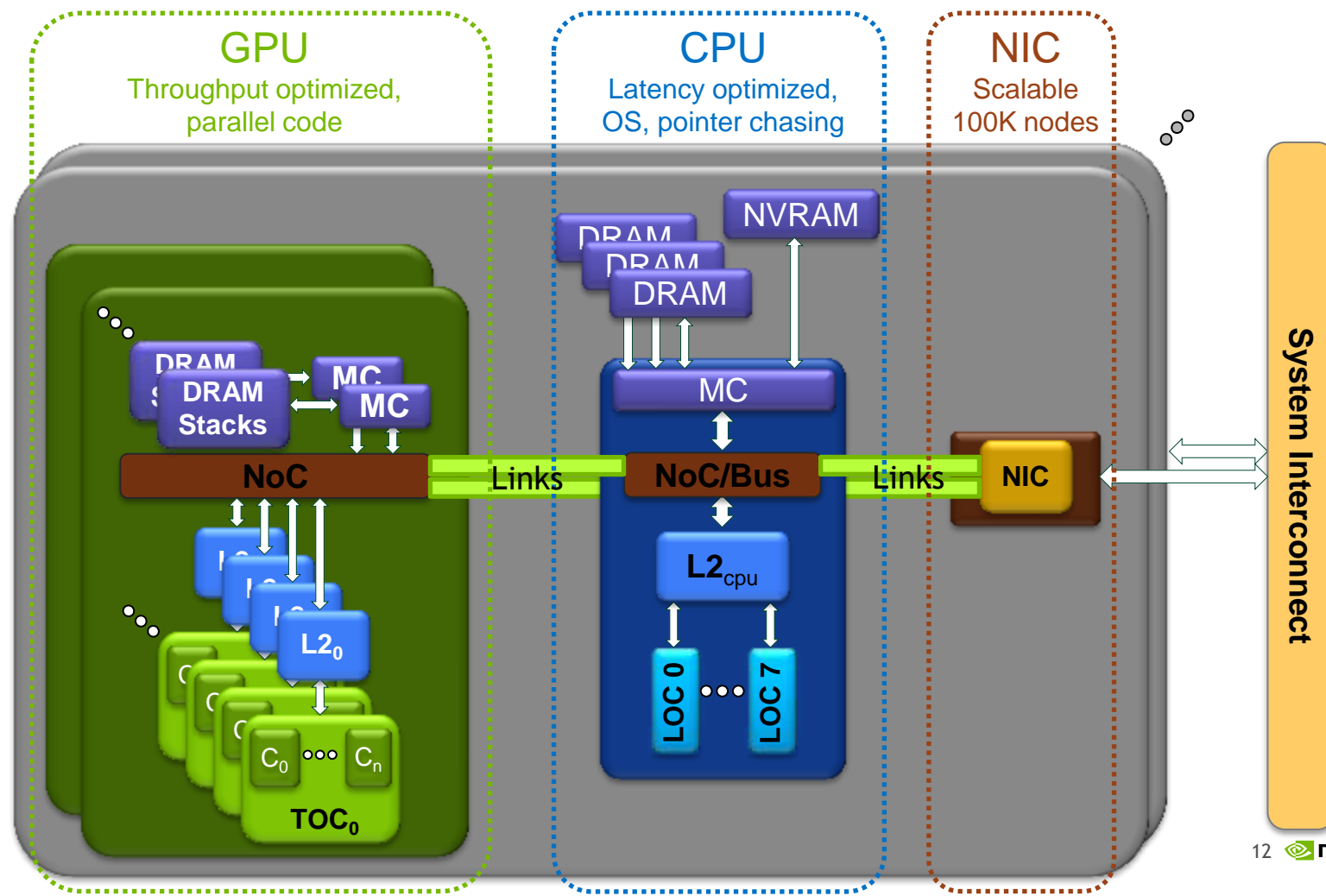Sequential Ops

**CPU**

11 NVIDIA

# Generic Future Node Model
## Three Building Blocks (GPU, CPU, Network)

**Direct Evolution**

- Programming Model Continuity

- Specialized Cores
  - GPU for parallel work
  - CPU for serial work

- Coherent memory system with Stacked, Bulk, & NVRAM

- Amortize non-parallel costs
  - Increase GPU:CPU
  - Smaller CPU

- *Can be one chip or MCM or multiple sockets*



**GPU**
Throughput optimized, parallel code

**CPU**
Latency optimized, OS, pointer chasing

**NIC**
Scalable 100K nodes

DRAM Stacks · DRAM Stacks · MC · MC · NoC · $L2_0$ · $C_0$ ··· $C_n$ · $TOC_0$

NVRAM · DRAM · DRAM · DRAM · MC · NoC/Bus · $L2_{cpu}$ · LOC 0 ··· LOC 7

Links · Links · NIC · System Interconnect

# 3 Ways to Program GPUs

**Applications**

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

**NVIDIA.**

# GPU DEVELOPER ECO-SYSTEM

### Numerical Packages

MATLAB
Mathematica
NI LabView
pyCUDA

### Debuggers & Profilers

cuda-gdb
NV Visual Profiler
Parallel Nsight
GPU Wizard
Allinea
TotalView

### GPU Compilers

C
C++
Fortran
Java
Python

### Auto-parallelizing & Cluster Tools

OpenACC
GPUDirect
OpenMP
Ocelot

### Libraries

BLAS
FFT
LAPACK
NPP
Video
Imaging
GPULib

### Consultants & Training

acceleware
GASS
STONE RIDGE TECHNOLOGY
ANEO
GPU Tech
TECH
ArrayFire
HPC
EM Photonics
SCALABLE GRAPHICS
WIPRO

### OEM Solution Providers

DELL
hp
IBM
CRAY
ASUS
SUPERMICRO
sgi
FUJITSU
BULL
APPRO
lenovo 联想
NEC

# DEVELOP ON GEFORCE, DEPLOY ON TESLA

**GeForce GPUs**

**Tesla K40/K80**

**Designed for Gamers & Developers**

**Designed for Cluster Deployment**

Available Everywhere

https://developer.nvidia.com/cuda-gpus

ECC
24x7 Runtime
GPU Monitoring
Cluster Management
GPUDirect-RDMA
Hyper-Q for MPI
3 Year Warranty
Integrated OEM Systems, Professional Support

# 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**

NVIDIA.