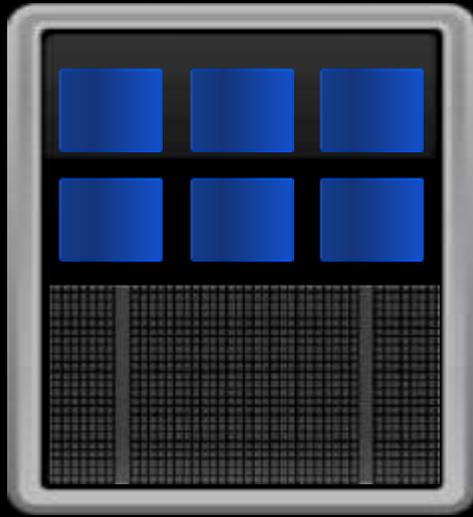


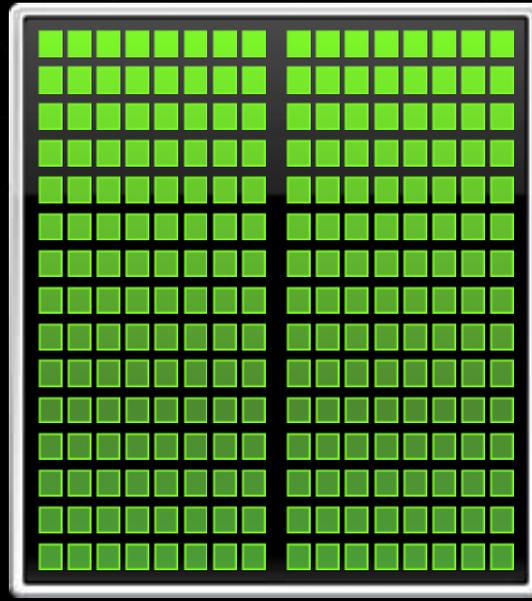
ACCELERATED COMPUTING

10X PERFORMANCE & 5X ENERGY EFFICIENCY

CPU
Optimized for
Serial Tasks



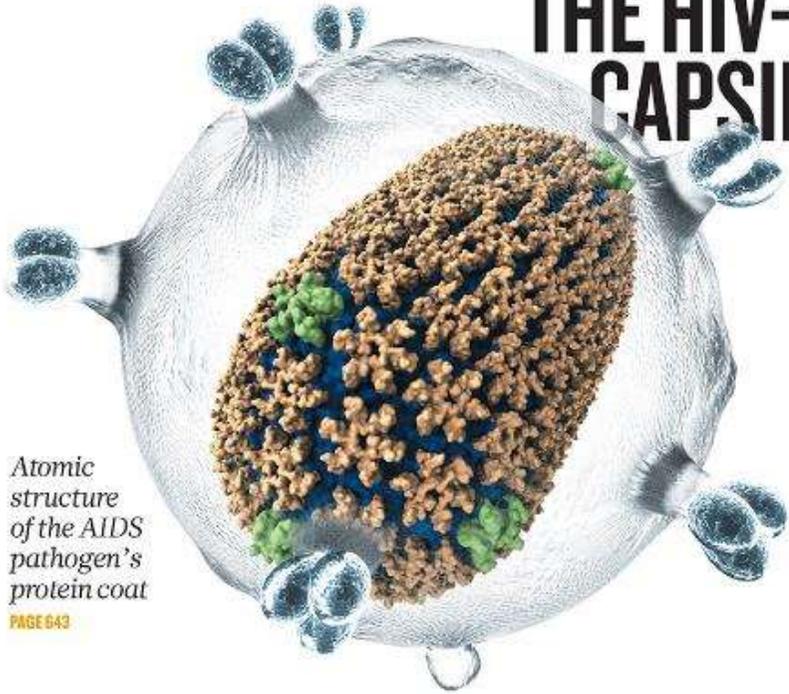
GPU Accelerator
Optimized for
Parallel Tasks



nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

THE HIV-1 CAPSID



Atomic
structure
of the AIDS
pathogen's
protein coat

PAGE 643

ACCELERATING DISCOVERIES

USING A SUPERCOMPUTER POWERED BY **3,000 TESLA PROCESSORS**, UNIVERSITY OF ILLINOIS SCIENTISTS PERFORMED THE FIRST ALL-ATOM SIMULATION OF THE HIV VIRUS AND DISCOVERED THE CHEMICAL STRUCTURE OF ITS CAPSID — “THE PERFECT TARGET FOR FIGHTING THE INFECTION.”

WITHOUT GPU, THE SUPERCOMPUTER WOULD NEED TO BE 5X LARGER FOR SIMILAR PERFORMANCE.

ACCELERATING INSIGHTS

“*Now You Can Build Google’s
\$1M Artificial Brain on the Cheap*”

WIRED

GOOGLE DATACENTER



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

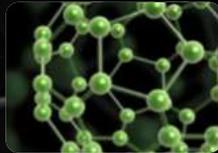
FROM HPC TO ENTERPRISE DATACENTERS



Oil & Gas



PETROBRAS



Higher Ed



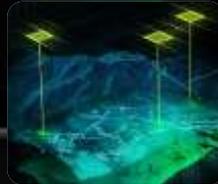
HARVARD
School of Engineering
and Applied Sciences



STANFORD
UNIVERSITY



UNIVERSITY OF
CAMBRIDGE



Government



Air Force
Research
Laboratory



Naval Research
Laboratory



Supercomputing



CSCS



NCSA



Tokyo Institute
of Technology



Lawrence Livermore
National Laboratory



Finance



Consumer Web





POPULAR GPU-ACCELERATED APPLICATIONS

- 02 Research: Higher Education and Supercomputing
 - COMPUTATIONAL CHEMISTRY AND BIOLOGY
 - BIOMEDICAL ANALYTICS
 - PHYSICS
 - WEATHER AND CLIMATE SIMULATION
- 06 Defense and Intelligence
- 07 Computational Finance
- 08 Manufacturing: CAD and CAE
 - COMPUTER AIDED DESIGN
 - COMPUTATIONAL FLUID DYNAMICS
 - COMPUTATIONAL STRUCTURAL MECHANICS
 - FLUID FLOW SIMULATION
- 10 Media and Entertainment
 - ANIMATION, RENDERING AND RENDERING
 - COLOR CORRECTION AND GRAF WORKSPACE
 - COMPOSITING, FINISHING AND EFFECTS
 - EDITING
 - ENCODING AND DIGITAL DISTRIBUTION
 - ON AIR GRAPHICS
 - ON SET, BEHIND THE SCENES AND EMULATOR
 - WEATHER SERVICES
- 14 Oil and Gas

Research: Higher Education and Supercomputing

COMPUTATIONAL CHEMISTRY AND BIOLOGY

Bioinformatics

Application	Description	Hardware Support	Operating System	Programming Language	Multi-Process	Availability
BerryCUDA	Sequence mapping software	Alignment of short sequencing reads	4-16x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 0.4.3
CDASH++	Open source software for Smith-Waterman protein database searches on GPUs	Parallel search of Smith-Waterman database	10-50x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 2.0.8
CUSMM	Parallelized short read aligner	Parallel, accurate long read aligner - gapped alignments to large genomes	16x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 1.0.0f
GPU-BLAST	Local search with fast k-mers	Prism algorithm according to k-mer, multi-processor	3-4x	T 2075, 2090, 410, 420, K20K	Single only	Available now Version 2.0.24
GPU-HMMER	Parallelized local and global search with profile Hidden Markov models	Parallel local and global search of Hidden Markov Models	40-100x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 2.0.3
mOUSA-NEMO	Ultrafast scalable motif discovery algorithm based on MEME	Scalable motif discovery algorithm based on MEME	4-10x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 3.0.11
SeqFlux	A GPU Accelerated Sequence Analysis Toolkit	Reference assembly, blast, motif-extraction, term, de-novo assembly	400x	T 2075, 2090, 410, 420, K20K	Yes	Available now
USENET	Open-source Smith-aligner for SOAPdenovo, suffix array based repeats filter and output	Fast short read alignment	9-10x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 1.11
WidML	Fits numerous linear models to a fixed design and response	Parallel linear regression on multiple arbitrary-shaped matrices	100x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 0.1-1

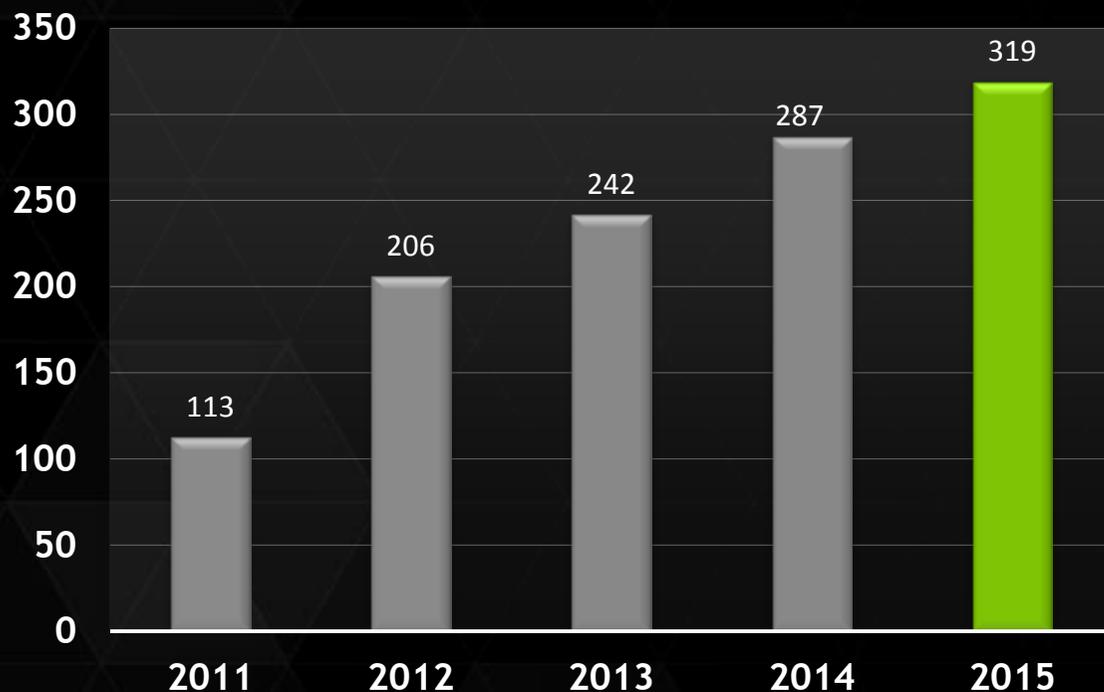
Molecular Dynamics

Application	Description	Hardware Support	Operating System	Programming Language	Multi-Process	Availability
Akane	Models molecular dynamics of biopolymers for simulations of proteins, DNA and lipids	Simulations on 100k GPUs	4-25x	T 2075, 2090, 410, 420, K20K	Single Only	Available now Version 1.0.40
ACMD	GPU simulation of molecular mechanics force fields, implicit and explicit solvent	Written for use on GPUs	140x (single GPU version only)	T 2075, 2090, 410, 420, K20K	Yes	Available now
AMBER	Suite of programs to simulate molecular dynamics on biomolecules	PHENIX, explicit and implicit solvent	80-140x (single GPU version)	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 12.0.1
DL-POLY	Simplest macromolecules, polymers, melt systems, etc on a distributed memory per-GPU computer	Two-body forces, Link-cell pairs, Rapid SPME forces, SHAKE-M	4x	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 0.1
CHARMM	MD package to simulate molecular dynamics on biomolecules	Implicit Sol, Explicit Sol, limited for OpenMP	700x	T 2075, 2090, 410, 420, K20K	Yes	In Development 0.12.2
CHARMM-Blue	Simulation of biomolecular molecules with complicated bond interactions	Implicit Sol, Explicit Sol, solvent	140x (single GPU version)	T 2075, 2090, 410, 420, K20K	Single only	Available now Version 6.8 in 0.12.2
HOOMD-blue	Particle dynamics package written primarily for GPUs	Written for GPUs	2x	T 2075, 2090, 410, 420, K20K	Yes	Available now
LAMMPS	Classical molecular dynamics package	Lammps-Jones, Morse, Buckingham, CHARMM, Tabulated, Custom green, ICH, Anisotropic, Gay-Born, BE-squared, Hybrid combination	3-10x	T 2075, 2090, 410, 420, K20K	Yes	Available now
MD	Designed for high performance simulation of large molecular systems	10k atom systems	4-16x (single GPU version)	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 2.0
OpenMM	Library and application for molecular dynamics for x86 with GPUs	Implicit and explicit solvent, nearest force	Implicit: 127-213x (single GPU version) Explicit: 16-10x (single GPU version)	T 2075, 2090, 410, 420, K20K	Yes	Available now Version 6.7.1

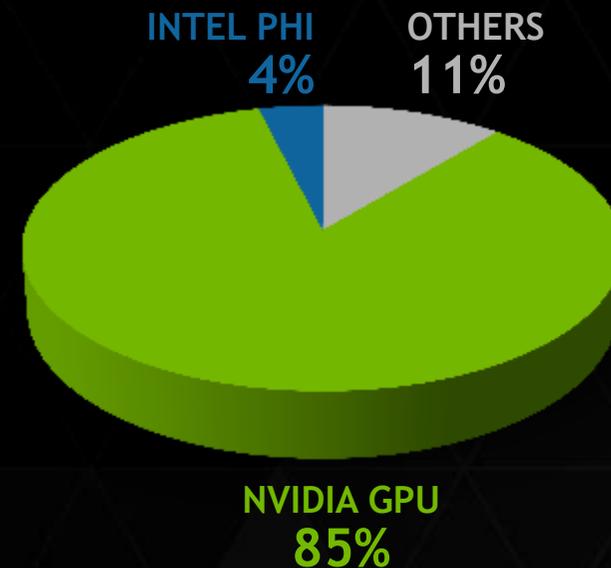
287 GPU-Accelerated Applications
www.nvidia.com/appscatalog

RAPID ADOPTION OF ACCELERATED COMPUTING

Hundreds of GPU Accelerated Apps



**NVIDIA GPU is
Accelerator of Choice**



HIGH GPU DENSITY SERVERS NOW MAINSTREAM



CRAY
THE SUPERCOMPUTER COMPANY

Cray CS-Storm
8 K80s per Node



Dell C4130
4 K80s per Node



HP SL270
8 K80s per Node



Quanta Computer

Quanta S2BV
4 K80s per Node

TESLA ACCELERATED COMPUTING PLATFORM

Data Center Infrastructure

System Solutions



Communication



Infrastructure Management



Development

Programming Languages



Development Tools



Software Solutions



GPU Accelerators
GPU Boost

Interconnect
GPU Direct
NVLink

System Management
NVML

Compiler Solutions
LLVM

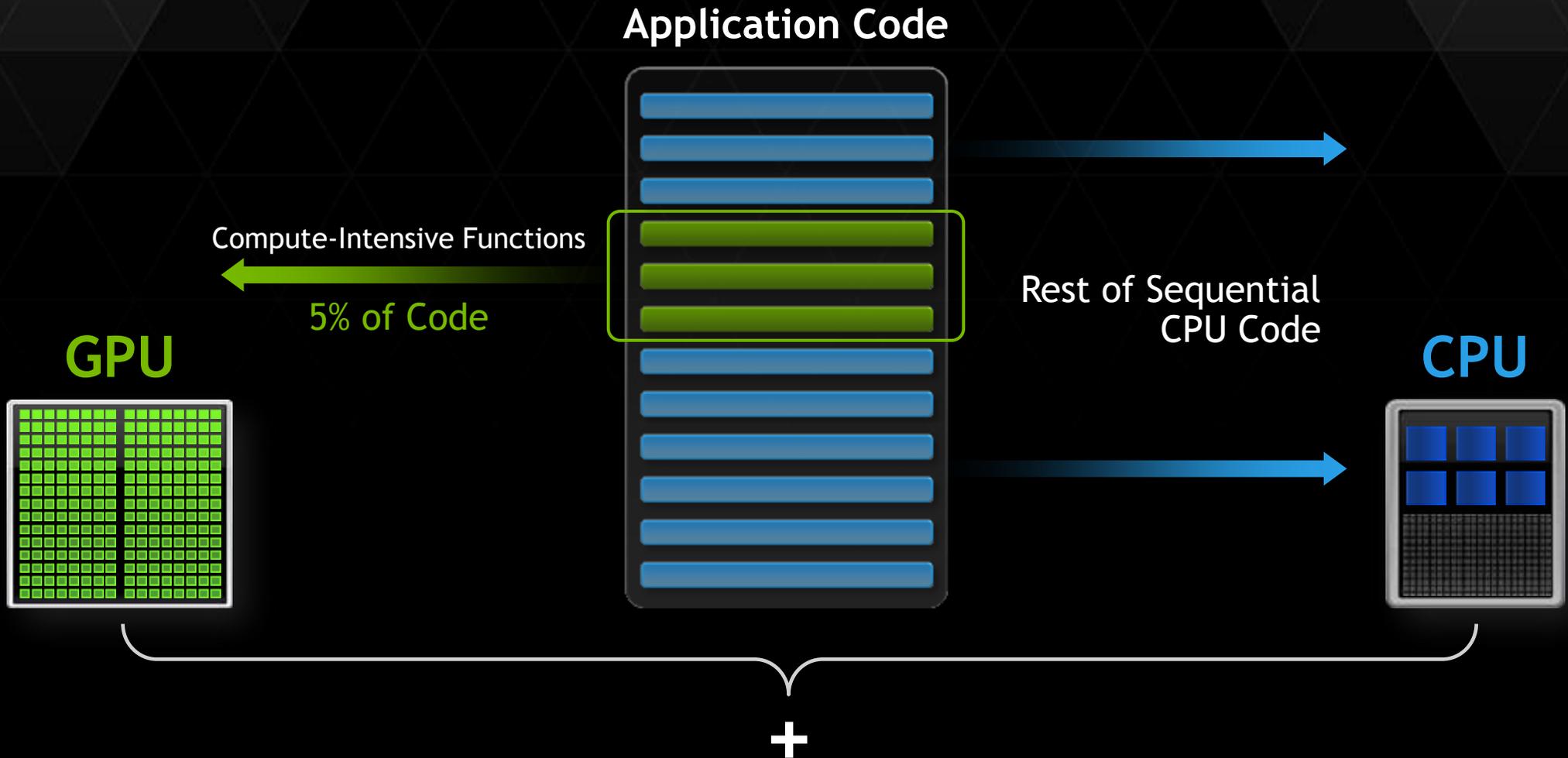
Profile and Debug
CUPTI

Accelerated Libraries
cuBLAS

Enterprise Services Support & Maintenance

TESLA PLATFORM FOR DEVELOPERS

HOW GPU ACCELERATION WORKS



3 WAYS TO PROGRAM GPUS

Applications

Libraries

“Drop-in”
Acceleration

OpenACC
Directives

Easily Accelerate
Applications

Programming
Languages

Maximum
Flexibility

GPU ACCELERATED LIBRARIES

“Drop-in” Acceleration for Your Applications

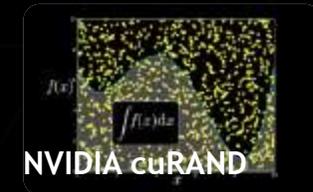
Linear Algebra

FFT, BLAS,
SPARSE, Matrix



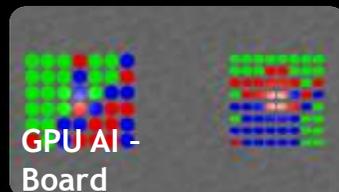
Numerical & Math

RAND, Statistics



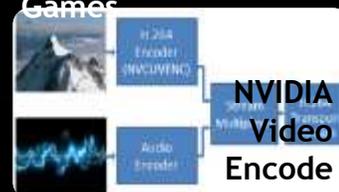
Data Struct. & AI

Sort, Scan, Zero Sum

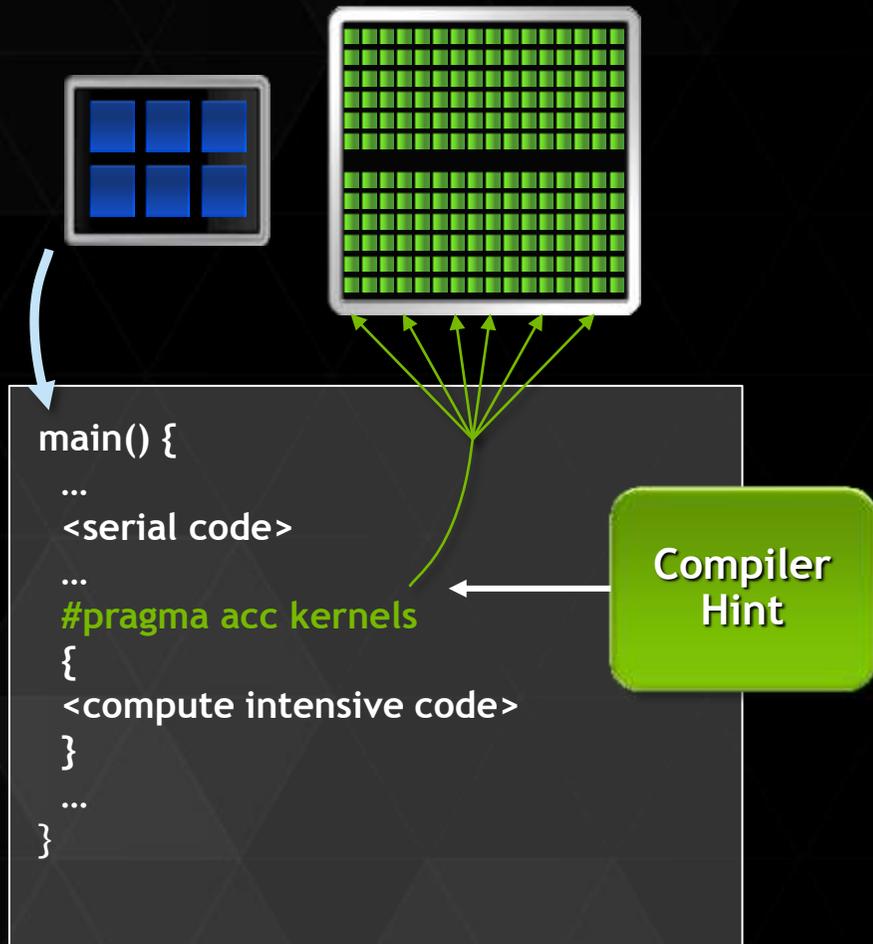


Visual Processing

Image & Video



OPENACC: OPEN, SIMPLE, PORTABLE



- Open Standard
- Easy, Compiler-Driven Approach
- Portable

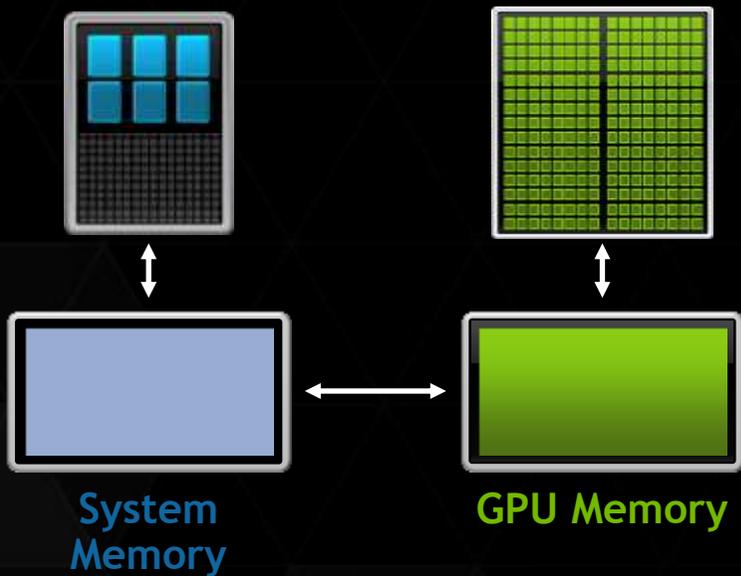
CAM-SE Climate
6x Faster on GPU
Top Kernel: 50% of Runtime



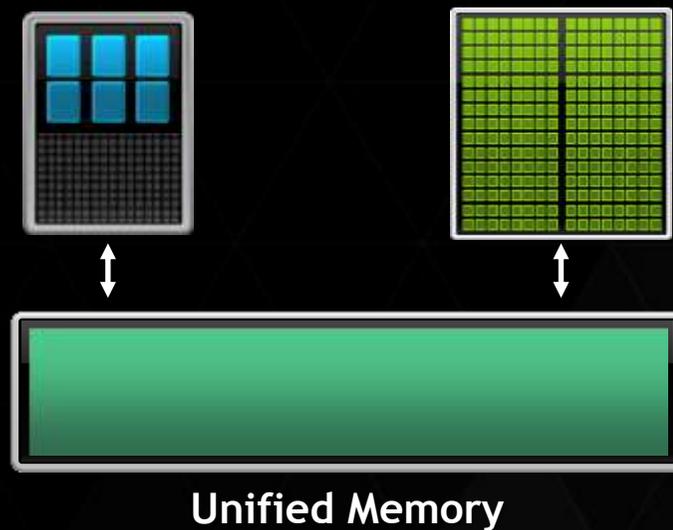
CUDA 6: UNIFIED MEMORY

Dramatically Lower Developer Effort

Developer View Today



Developer View With Unified Memory



SUPER SIMPLIFIED MEMORY MANAGEMENT CODE

CPU Code

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    data = (char *)malloc(N);  
  
    fread(data, 1, N, fp);  
  
    qsort(data, N, 1, compare);  
  
    use_data(data);  
  
    free(data);  
}
```

CUDA 6 Code with Unified Memory

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    cudaMallocManaged(&data, N);  
  
    fread(data, 1, N, fp);  
  
    qsort<<<...>>>(data, N, 1, compare);  
    cudaDeviceSynchronize();  
  
    use_data(data);  
  
    cudaFree(data);  
}
```

COMMON PROGRAMMING MODELS ACROSS MULTIPLE CPUS

Libraries



Compiler Directives

OpenACC



Programming Languages



GPU DEVELOPER ECO-SYSTEM

Numerical Packages

MATLAB
Mathematica
NI LabView
pyCUDA

Debuggers & Profilers

cuda-gdb
NV Visual Profiler
Parallel Nsight
Visual Studio
Allinea
TotalView

GPU Compilers

C
C++
Fortran
Java
Python

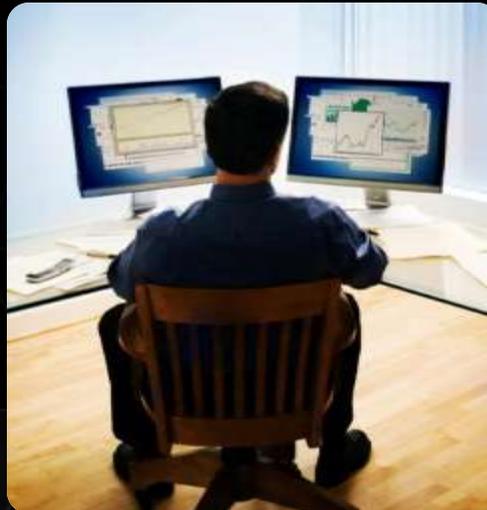
Auto-parallelizing & Cluster Tools

OpenACC
mCUDA
OpenMP
Ocelot

Libraries

BLAS
FFT
LAPACK
NPP
Video
Imaging
GPULib

Consultants & Training



OEM Solution Providers



DEVELOP ON GEFORCE, DEPLOY ON TESLA

GeForce GPUs

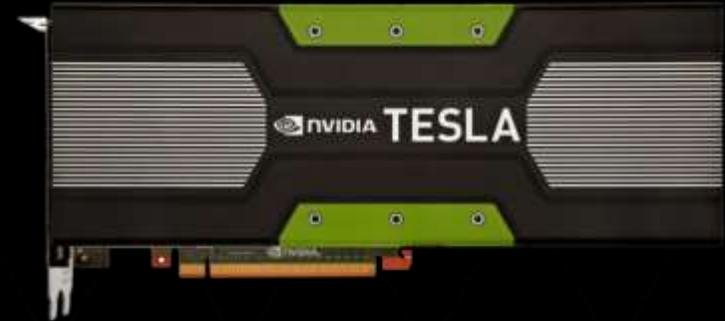


Designed for Gamers & Developers

Available Everywhere

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

Tesla K40/K80



Designed for Cluster Deployment

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

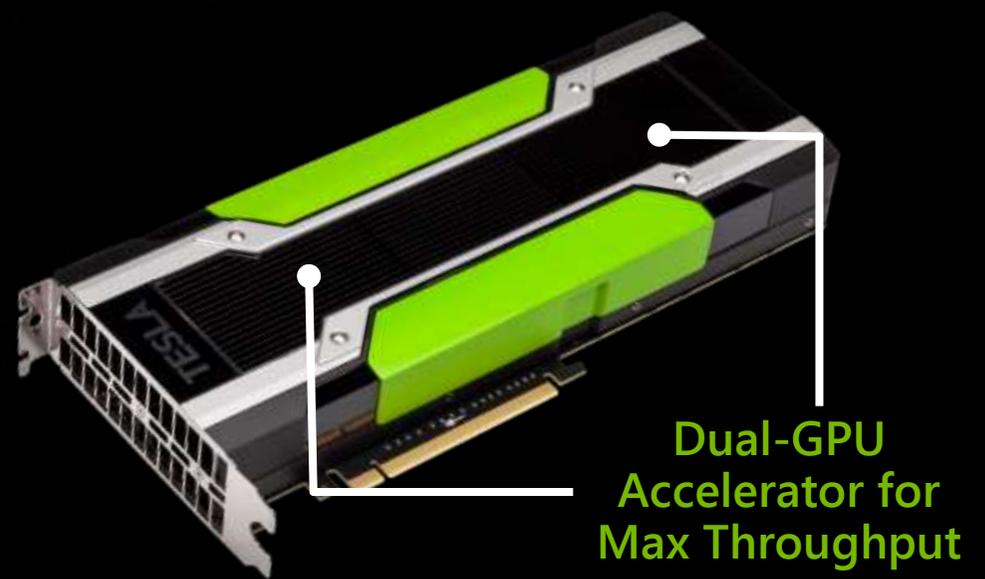
700+ University Courses
In **62** Countries



ACCELERATED COMPUTING ROADMAP

TESLA K80

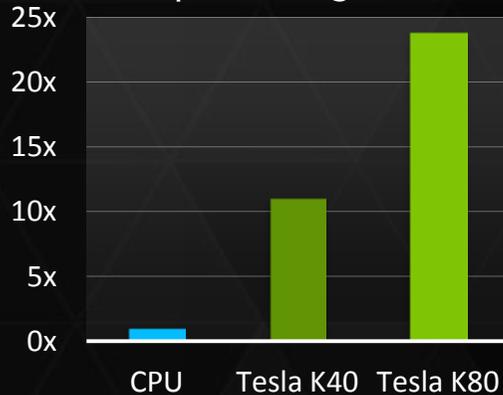
WORLD'S FASTEST ACCELERATOR
FOR DATA ANALYTICS AND
SCIENTIFIC COMPUTING



2x Faster

2.9 TF | 4992 Cores | 480 GB/s

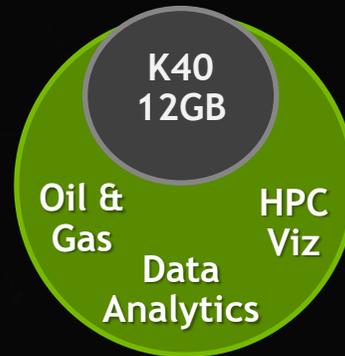
Deep Learning: Caffe



Double the Memory

Designed for Big Data Apps

24GB



Maximum Performance

Dynamically Maximize Perf for
Every Application

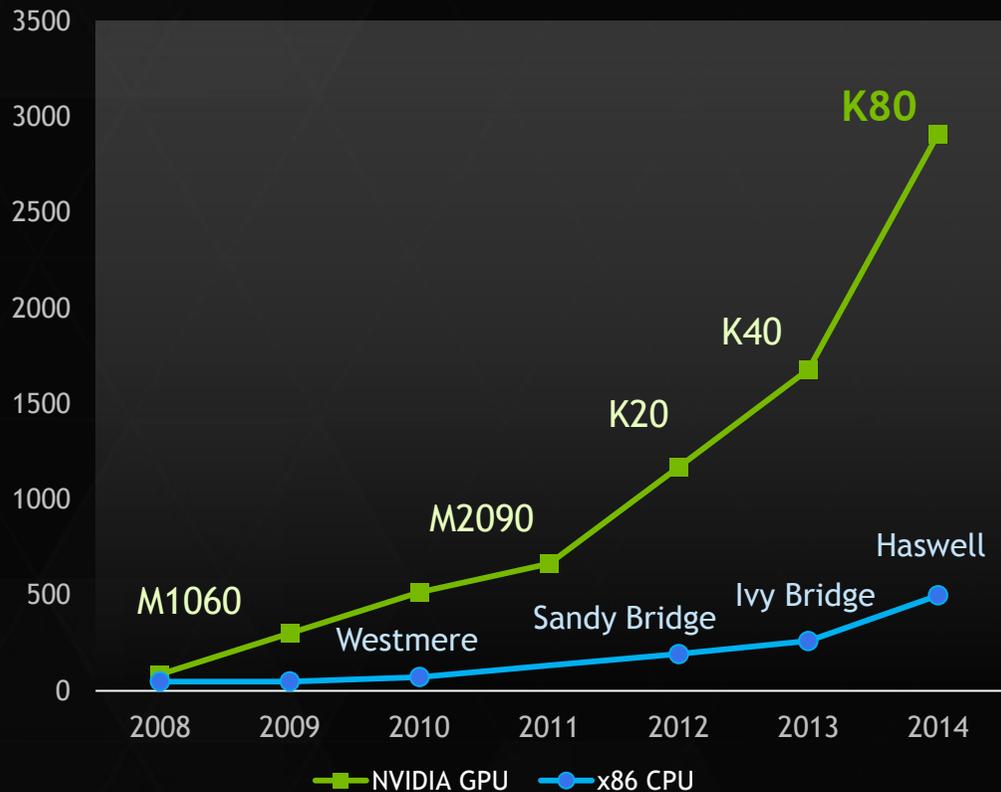


GPU Boost

PERFORMANCE LEAD CONTINUES TO GROW

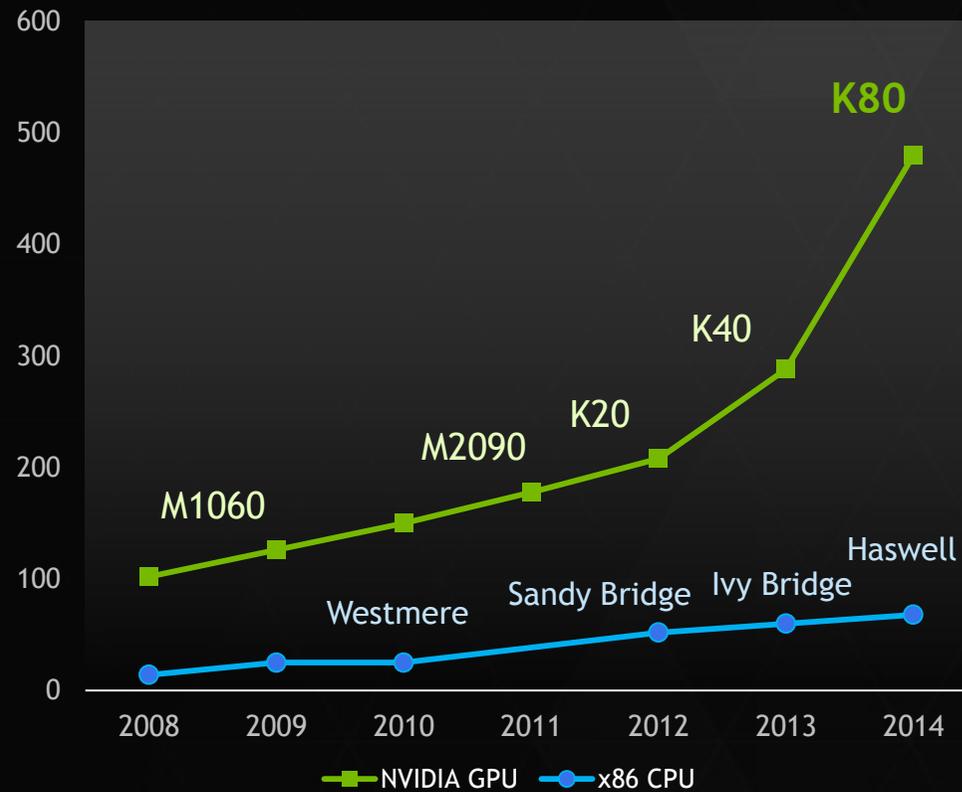
Peak Double Precision FLOPS

GFLOPS

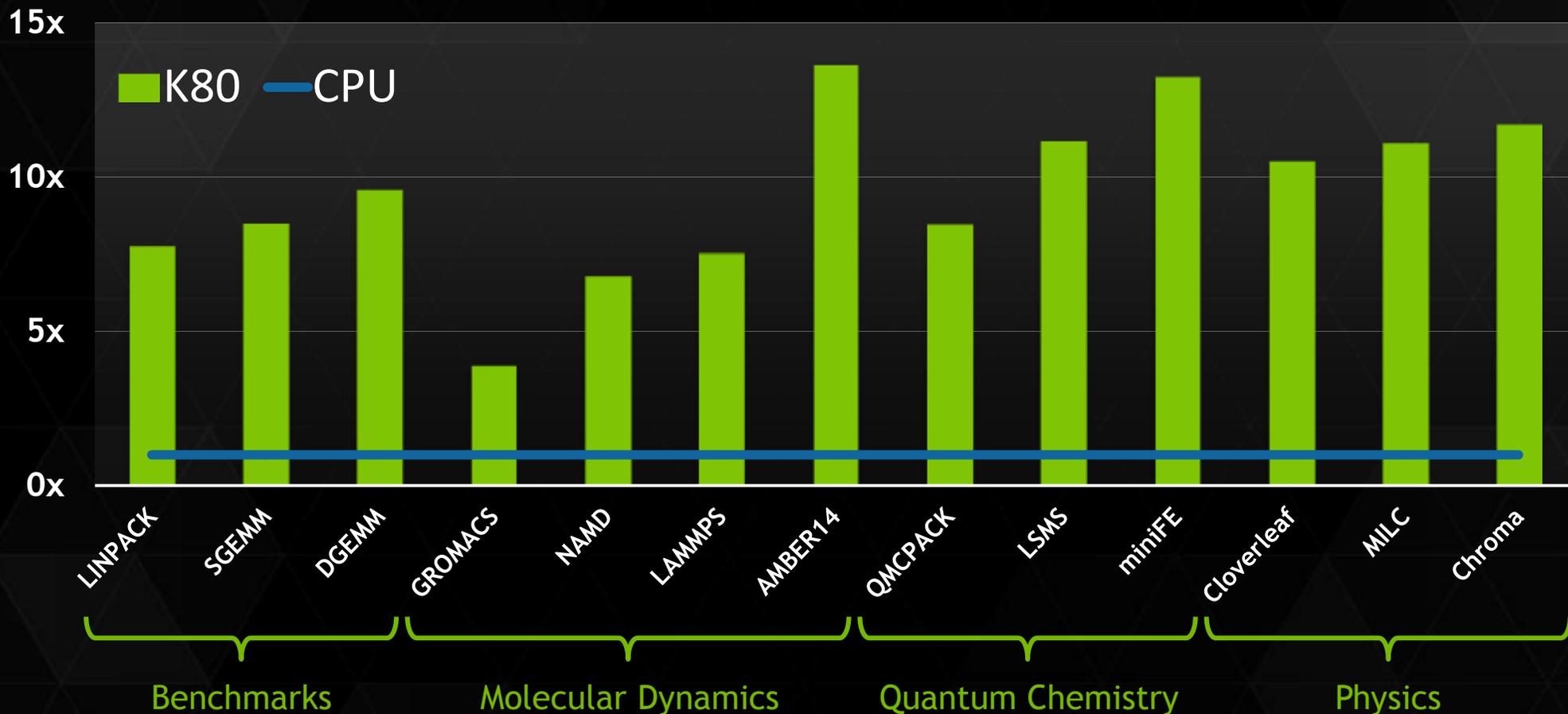


Peak Memory Bandwidth

GB/s

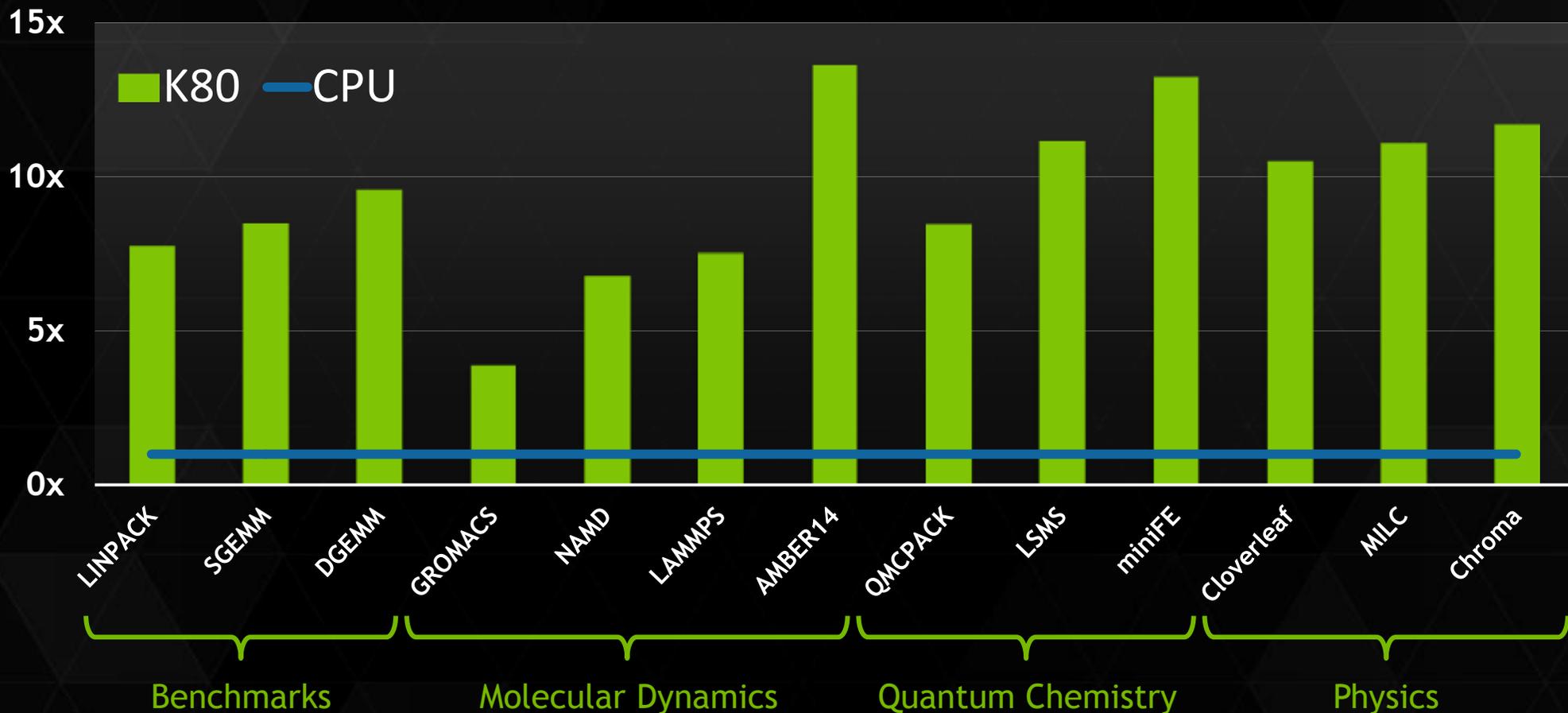


TESLA K80: 10X FASTER ON SCIENTIFIC APPS



CPU: 12 cores, E5-2697v2 @ 2.70GHz. 64GB System Memory, CentOS 6.2
GPU: Single Tesla K80, Boost enabled

TESLA K80: 10X FASTER ON REAL-WORLD APPS



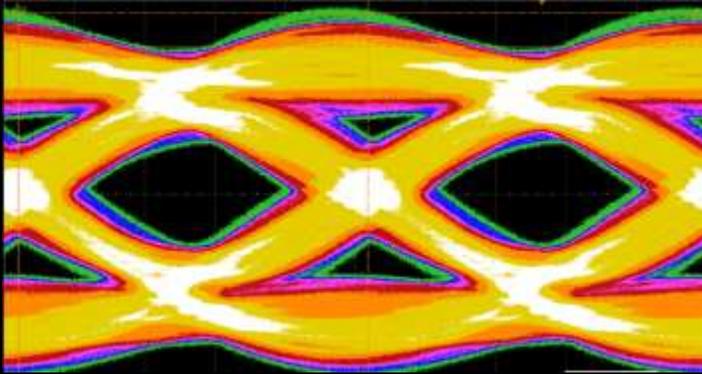
CPU: 12 cores, E5-2697v2 @ 2.70GHz. 64GB System Memory, CentOS 6.2
GPU: Single Tesla K80, Boost enabled

GPU ROADMAP



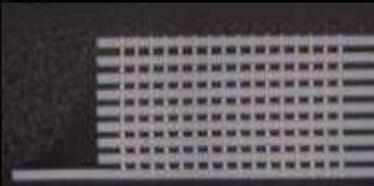
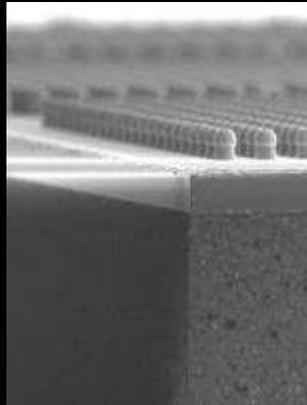
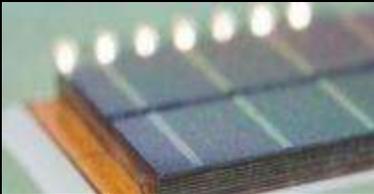
PASCAL GPU FEATURES

NVLINK AND STACKED MEMORY



NVLINK

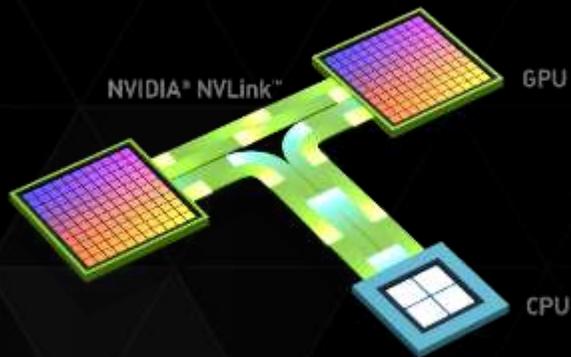
- GPU high speed interconnect
- 80-200 GB/s



3D Stacked Memory

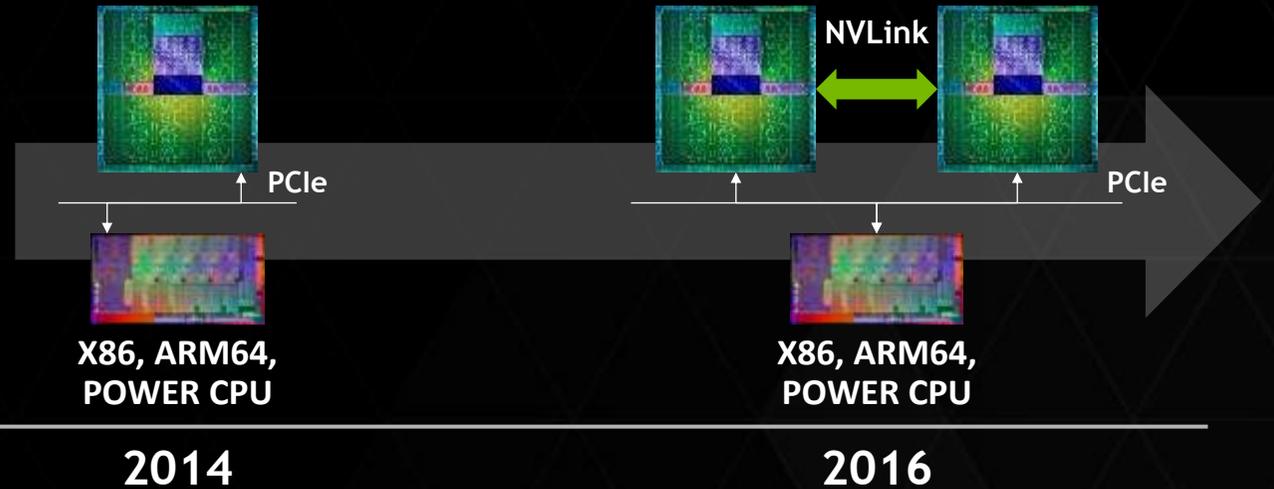
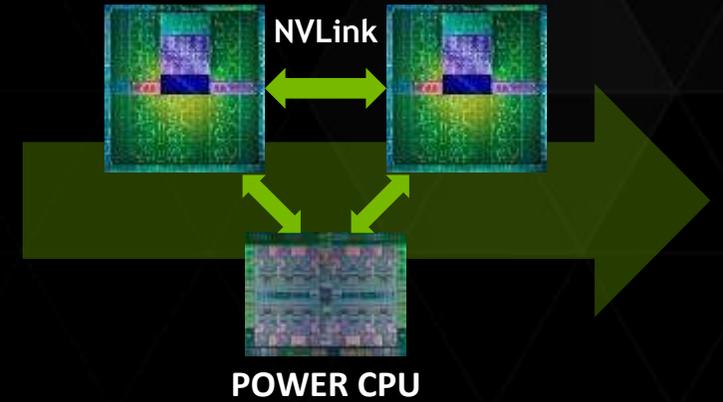
- 4x Higher Bandwidth (~1 TB/s)
- 3x Larger Capacity
- 4x More Energy Efficient per bit

NVLINK HIGH-SPEED GPU INTERCONNECT



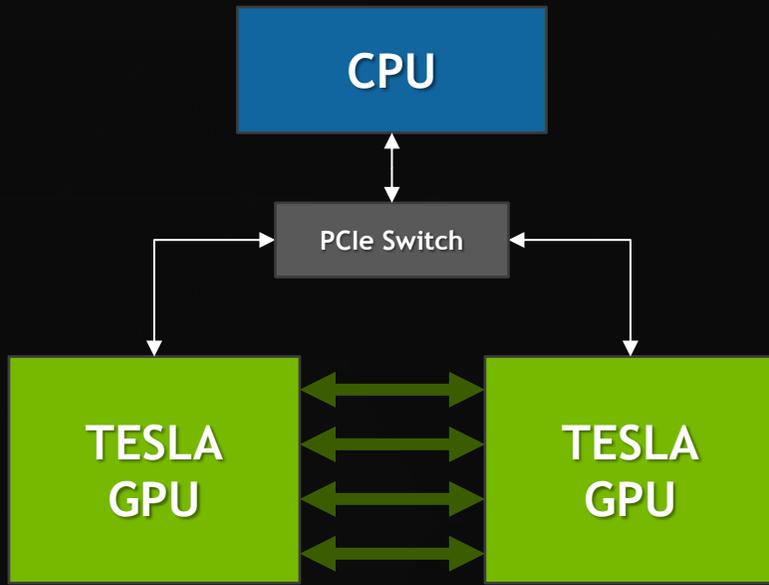
KEPLER GPU

PASCAL GPU



NVLINK UNLEASHES MULTI-GPU PERFORMANCE

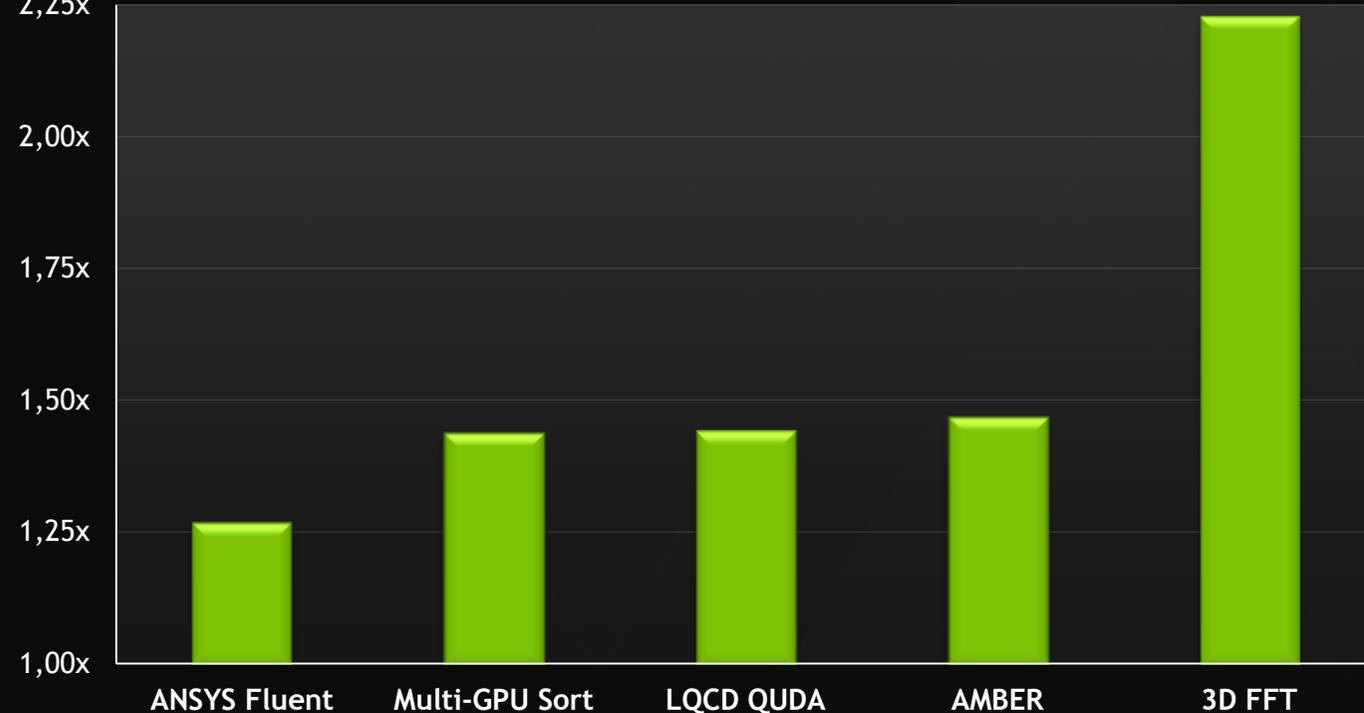
GPUs Interconnected with NVLink



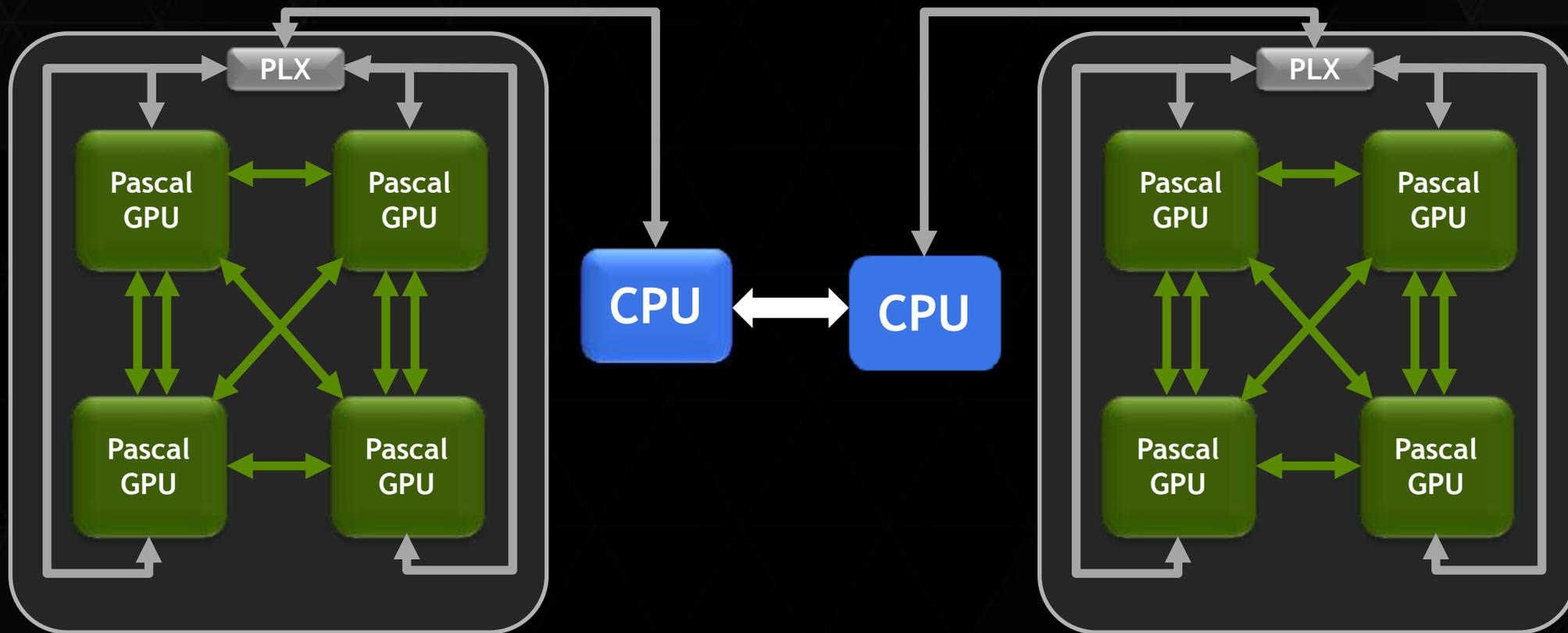
5x Faster than
PCIe Gen3 x16

Over 2x Application Performance Speedup When Next-Gen GPUs Connect via NVLink Versus PCIe

Speedup vs
PCIe based Server
2,25x



EXAMPLE: 8-GPU SERVER WITH NVLINK



↔ NVLINK 20GB/s
↔ PCIe x16 Gen 3

US TO BUILD TWO FLAGSHIP SUPERCOMPUTERS POWERED BY THE TESLA PLATFORM



100-300 PFLOPS Peak

10x in Scientific App Performance

IBM POWER9 CPU + NVIDIA Volta GPU

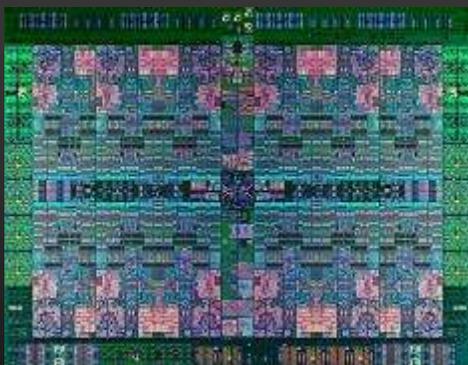
NVLink High Speed Interconnect

40 TFLOPS per Node, >3,400 Nodes

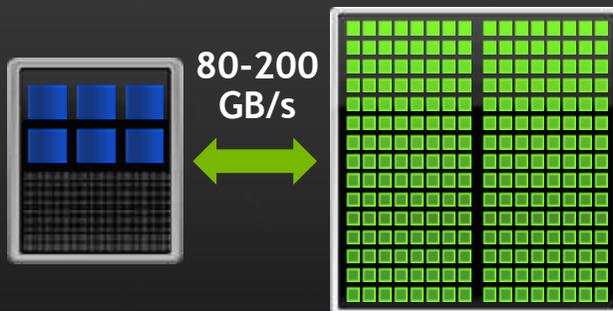
2017

Major Step Forward on the Path to Exascale

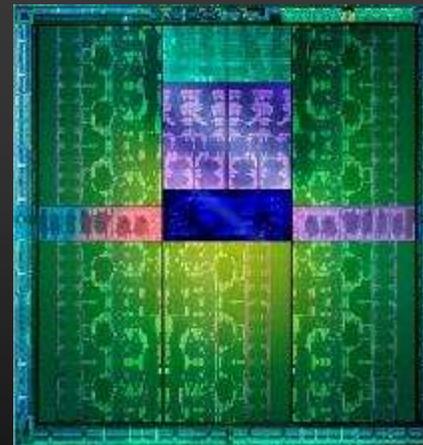
ACCELERATED COMPUTING 5X HIGHER ENERGY EFFICIENCY



IBM POWER CPU
Most Powerful Serial Processor

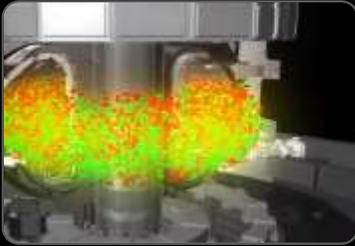


NVIDIA NVLink
Fastest CPU-GPU Interconnect



NVIDIA Volta GPU
Most Powerful Parallel Processor

CORAL: BUILT FOR GRAND SCIENTIFIC CHALLENGES



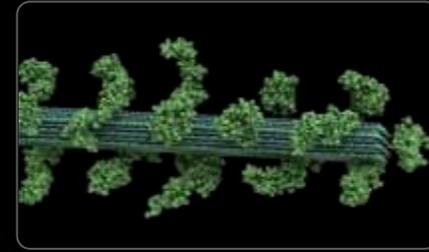
Fusion Energy

Role of material disorder, statistics, and fluctuations in nanoscale materials and systems.



Climate Change

Study climate change adaptation and mitigation scenarios; realistically represent detailed features

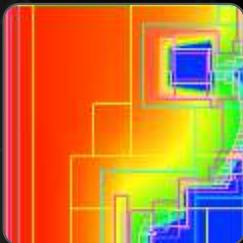


Biofuels

Search for renewable and more efficient energy sources

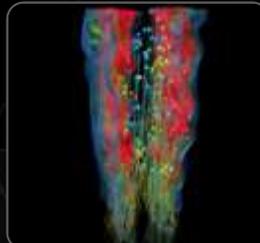
Astrophysics

Radiation transport – critical to astrophysics, laser fusion, atmospheric dynamics, and medical imaging



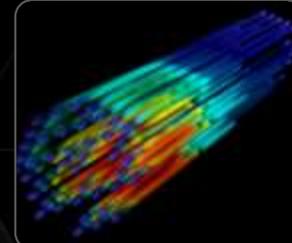
Combustion

Combustion simulations to enable the next gen diesel/bio-fuels to burn more efficiently



Nuclear Energy

Unprecedented high-fidelity radiation transport calculations for nuclear energy applications



IN-SITU VISUALIZATION

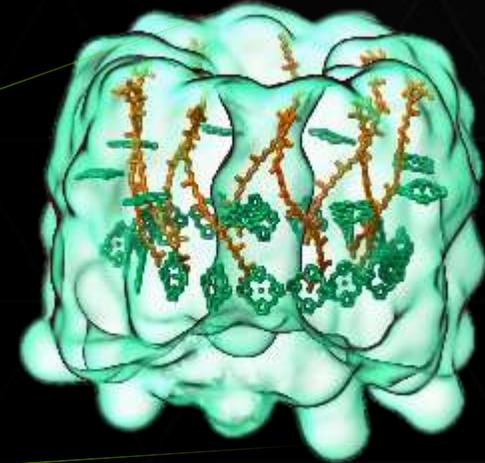
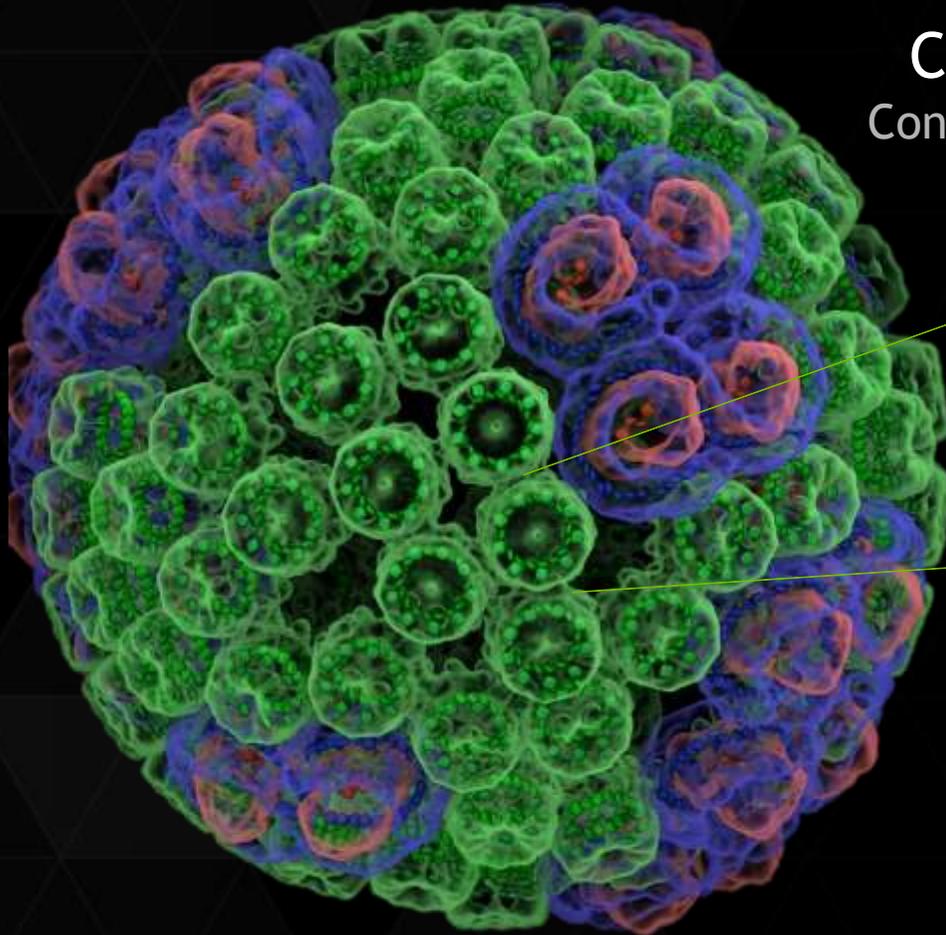
Enabling Visualization with the Tesla Platform



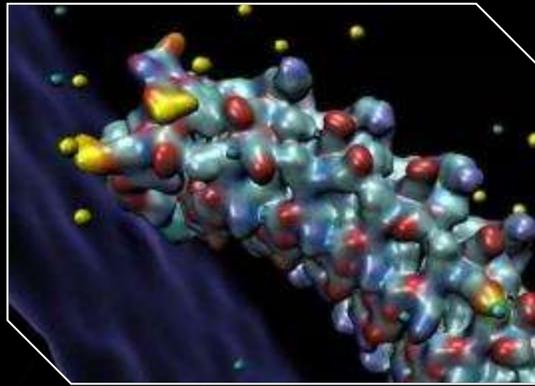
Theoretical and Computational Biophysics Group
University of Illinois at Urbana-Champaign

Chromatophore

Converts Light to Energy



WORLD'S LARGEST IN-SITU HPC VISUALIZATION



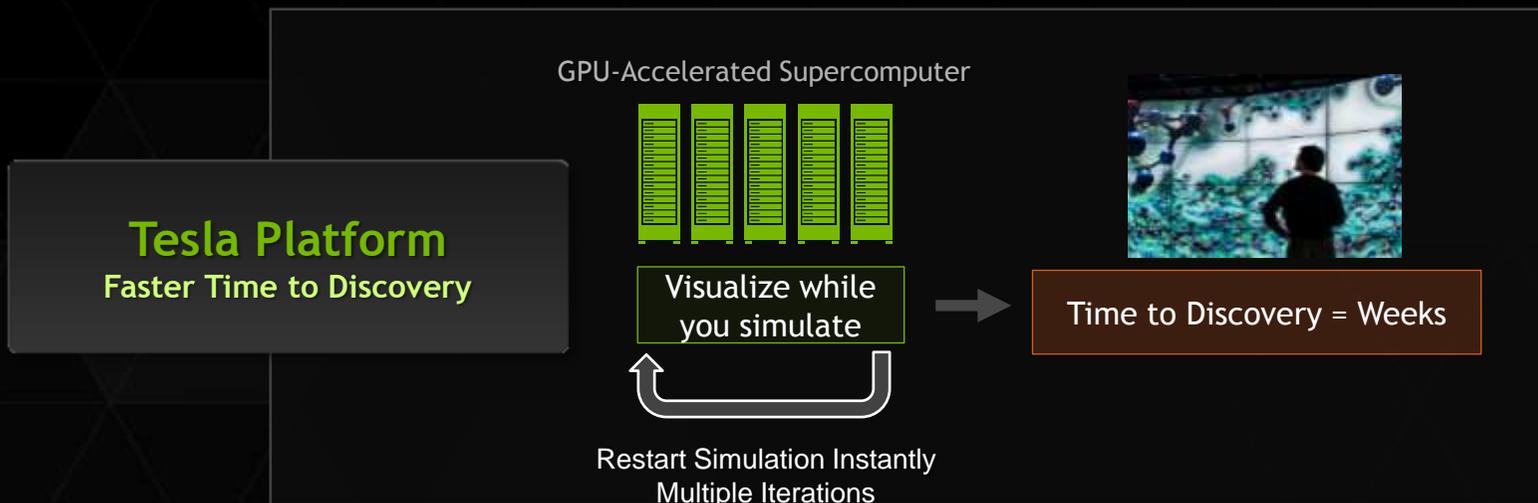
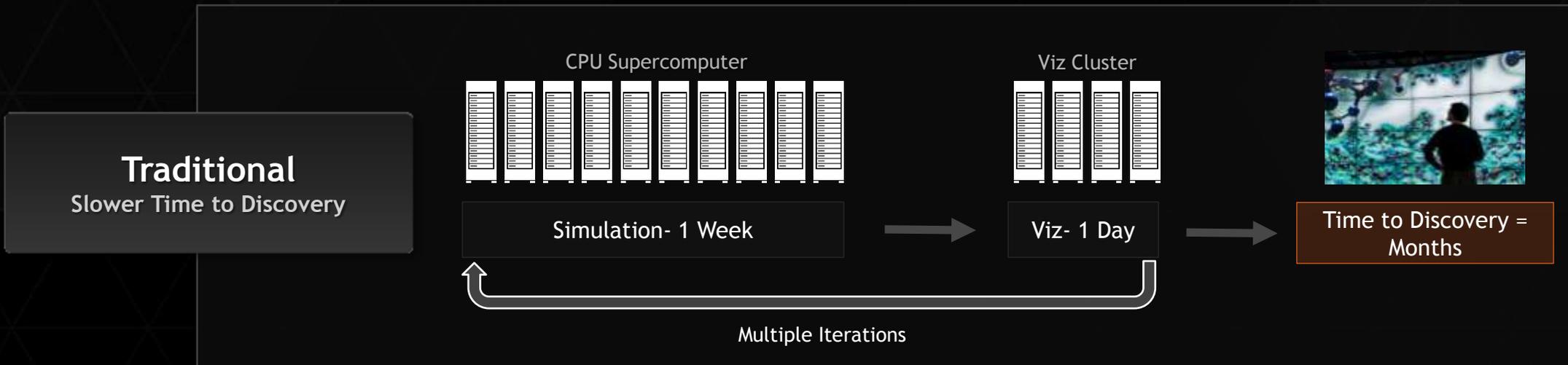
2048 GPU Nodes on CSCS Piz Daint

Galaxy formation and
Molecular Dynamics

Simulation + Visualization



VISUALIZE DATA INSTANTLY FOR FASTER SCIENCE

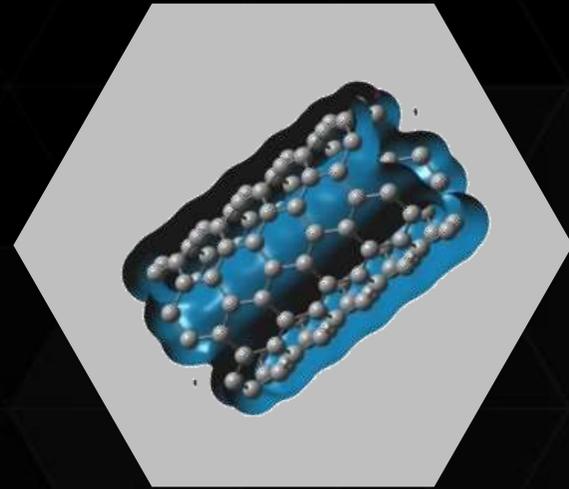
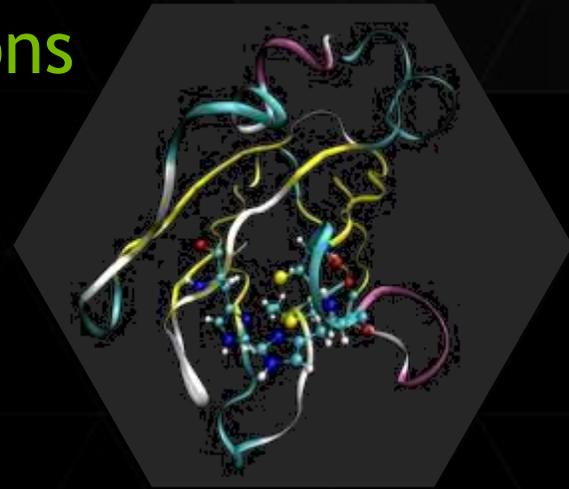


APPLICATIONS & CUSTOMER SUCCESSES

LIFE & MATERIAL SCIENCES

Overview of Accelerated Applications

- ▶ MD: All key codes are GPU-accelerated
 - ▶ **ACEMD***, **AMBER (PMEMD)***, BAND, CHARMM, DESMOND, ESPRESSO, Folding@Home, GPUgrid.net, GROMACS, HALMD, **HOOMD-Blue***, LAMMPS, Lattice Microbes, mdcore, NAMD, OpenMM, SOP-GPU
 - ▶ Great multi-GPU performance!
 - ▶ Focus: on dense (up to 16) GPU nodes & large # of GPU nodes
- ▶ QC: All key codes are ported or optimizing:
 - ▶ GPU-accelerated and available today:
 - ▶ ABINIT, ACES III, ADF, BigDFT, CP2K, GAMESS, Quantum Espresso/PWscf, MOLCAS, MOPAC2012, NWChem, QUICK, Q-Chem, **TeraChem***
 - ▶ Active GPU acceleration projects:
 - ▶ CASTEP, CPMD, GAMESS, Gaussian, NWChem, ONETEP, Quantum Supercharger Library, VASP & more
 - ▶ Focus: on using GPU-accelerated math libraries, OpenACC directives



REVOLUTIONIZING SCIENTIFIC COMPUTING

AMBER Molecular Dynamics Simulation
DHFR NVE Benchmark



64 Sandy Bridge CPUs
58 ns/day



Server with 2 Tesla K80
220 ns/day

ACCELERATING SIGNAL & VIDEO ANALYTICS

Real-time HD video
enhancements and analytics

Made possible only with GPUs



Video surveillance with faster
than real time analytics

12x faster with GPUs



Unmanned submarine with
accelerated sonar processing

50-100x speed up over CPU



Faster satellite image processing
for actionable intelligence

12x faster using GPUs

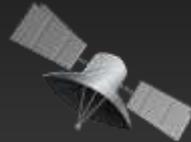


MISSION PLANNING WITH REAL-TIME LINE OF SIGHT

Video Data



Image Data



Signal Data

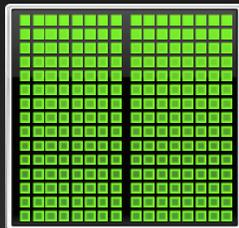


CPU



1 Computation/Second
Delayed Response

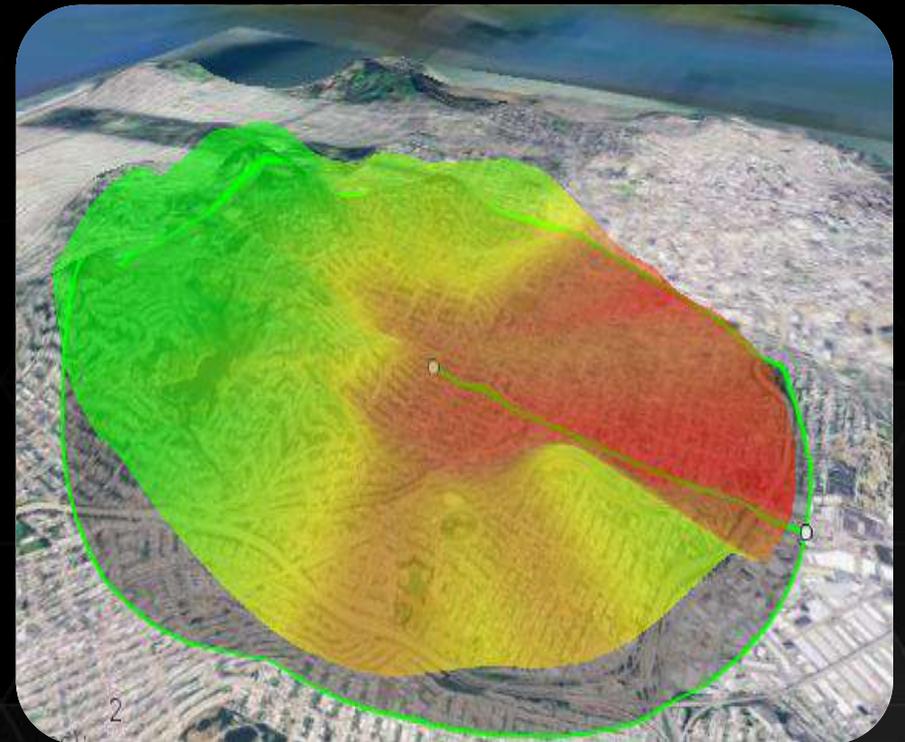
GPU



100 Computations/Second
Real-Time Response

LUCIAD

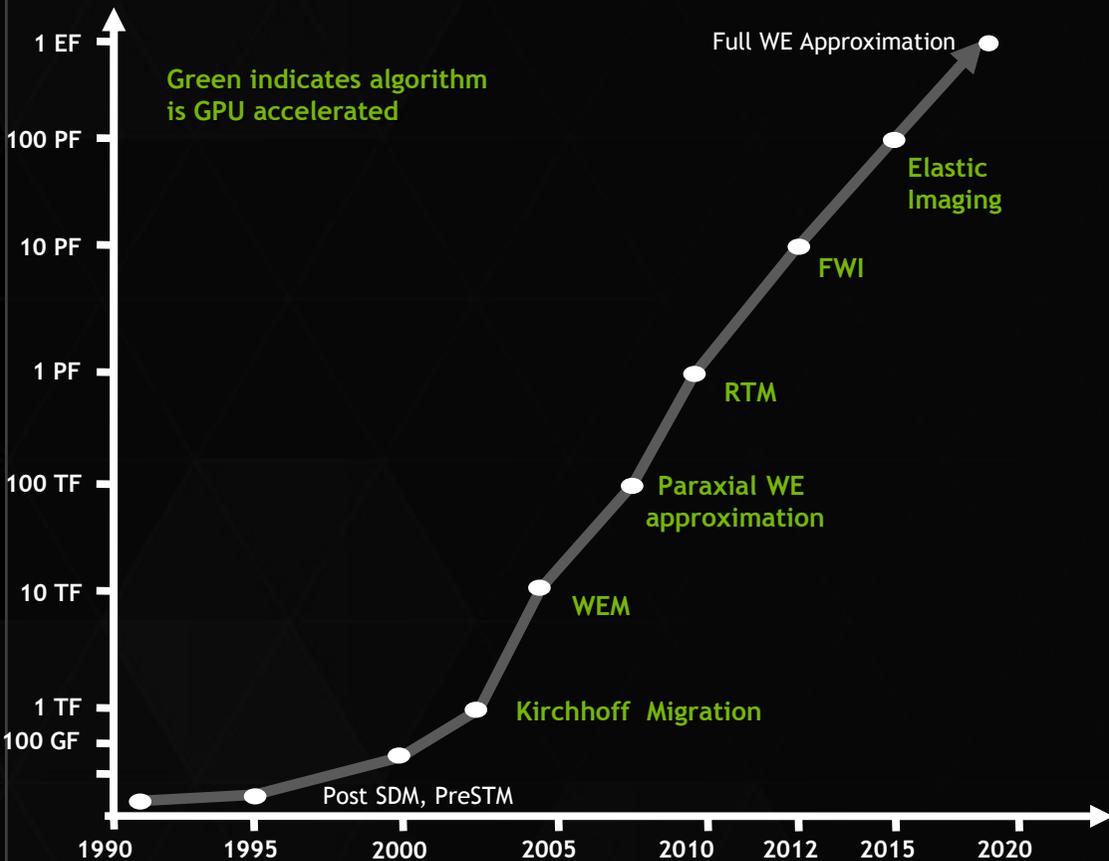
World Leader in Geospatial Situational Awareness



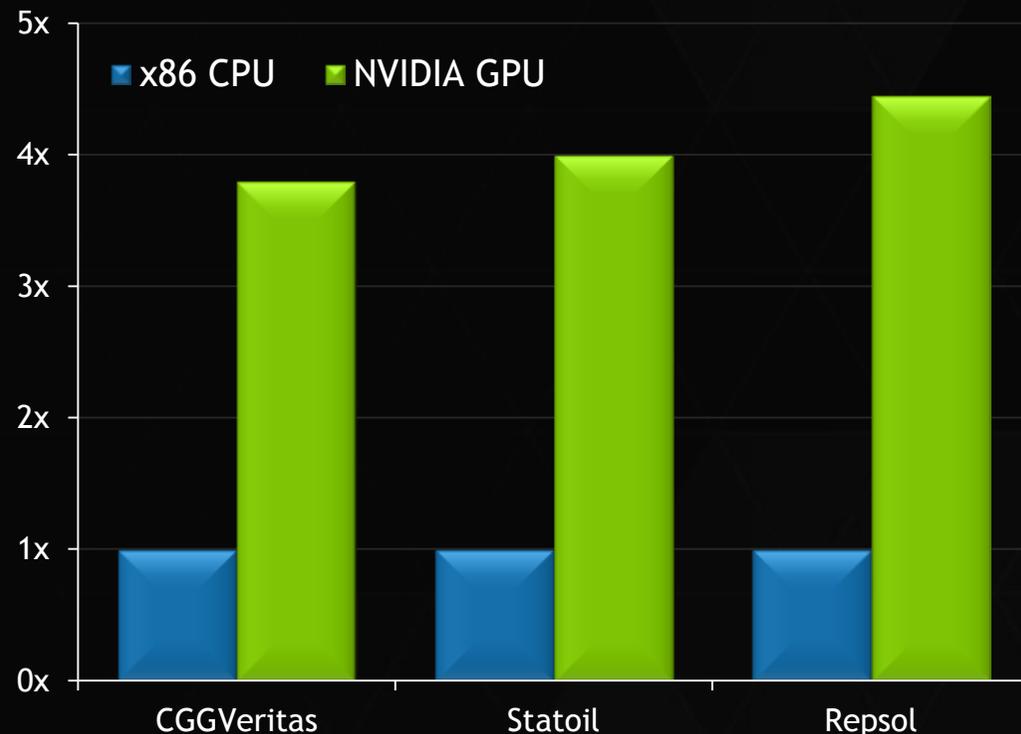
<http://www.luciad.com/>

KEY O&G APPS ACCELERATED ON GPUS

Most Workloads Ported to GPUs



Higher Throughput with GPUs



CGGVeritas: 2013 Rice Oil & Gas HPC Workshop
Statoil: [Maximizing TTI RTM Throughput for CPU + GPU](#); EAGE Conference (2013)
Repsol: (Test case: TT-A) [NVIDIA GTC 2014](#)

WORLD'S FASTEST ENTERPRISE SUPERCOMPUTER

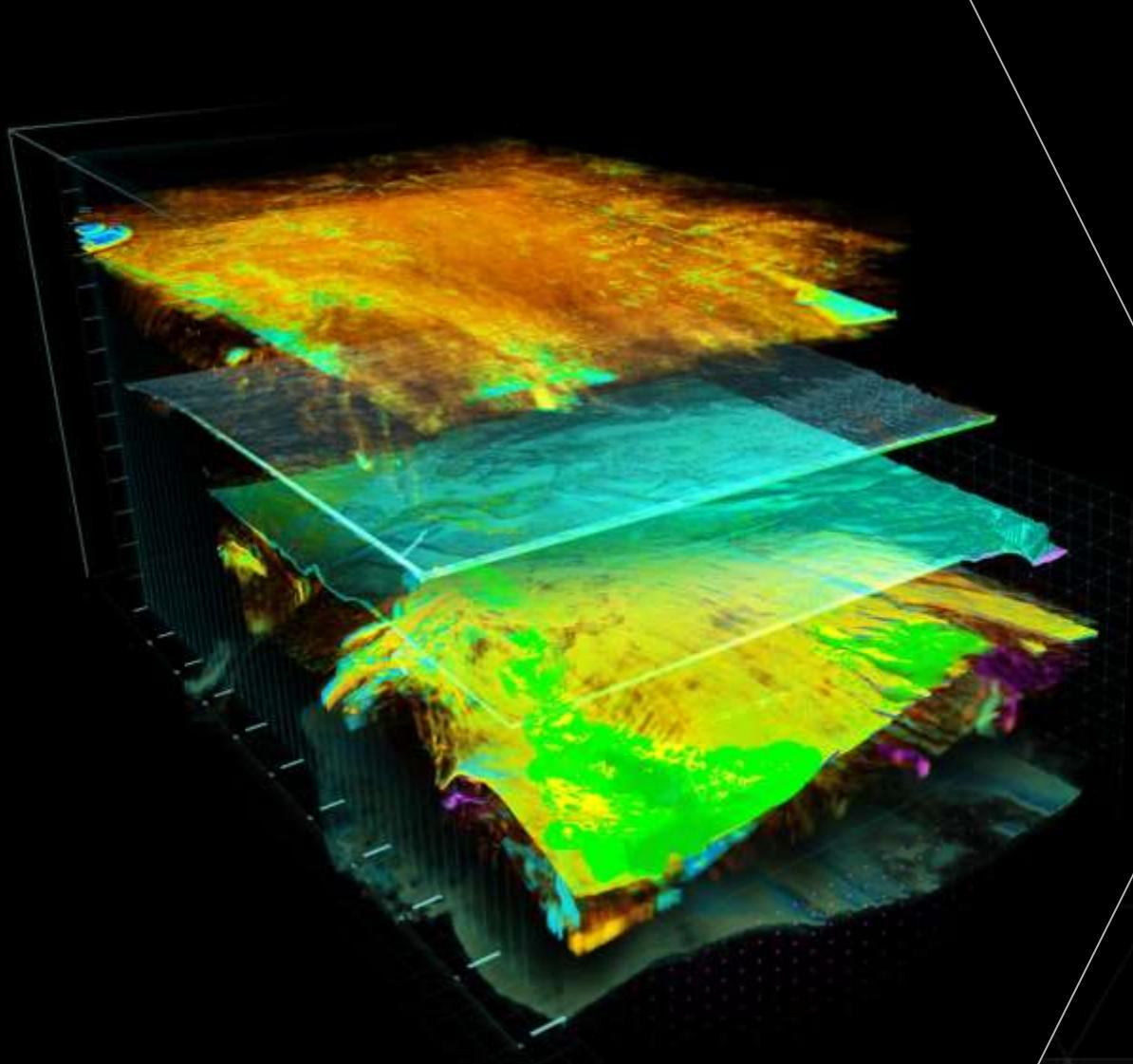


3 Petaflops Linpack Performance

Most Energy-Efficient Petascale
System in the World

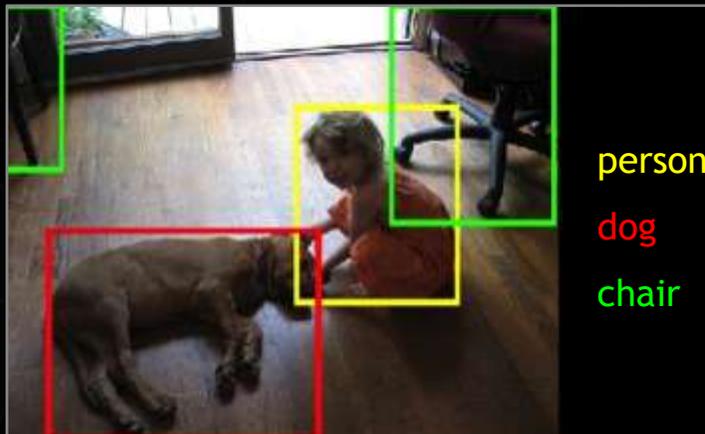
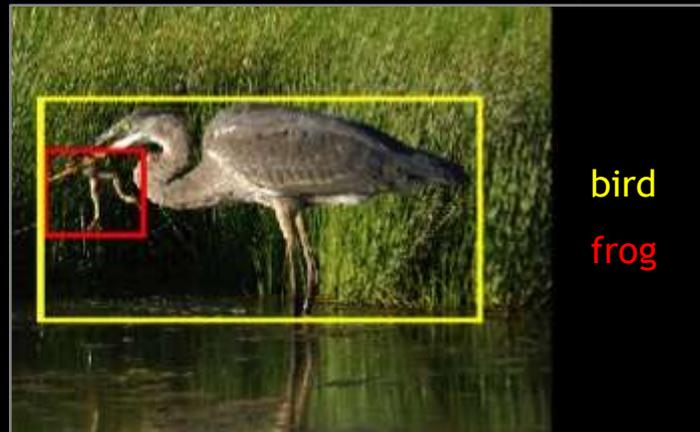
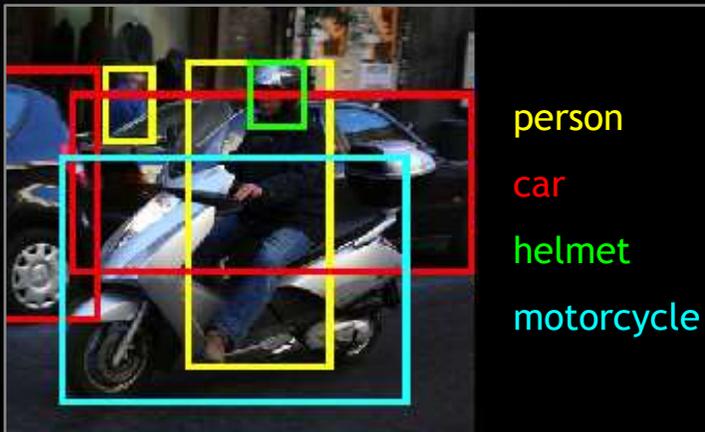
3,000 NVIDIA Tesla K20X GPU Accelerators

Maximizing Opportunity for Oil Discovery with
GPU-powered Supercomputer

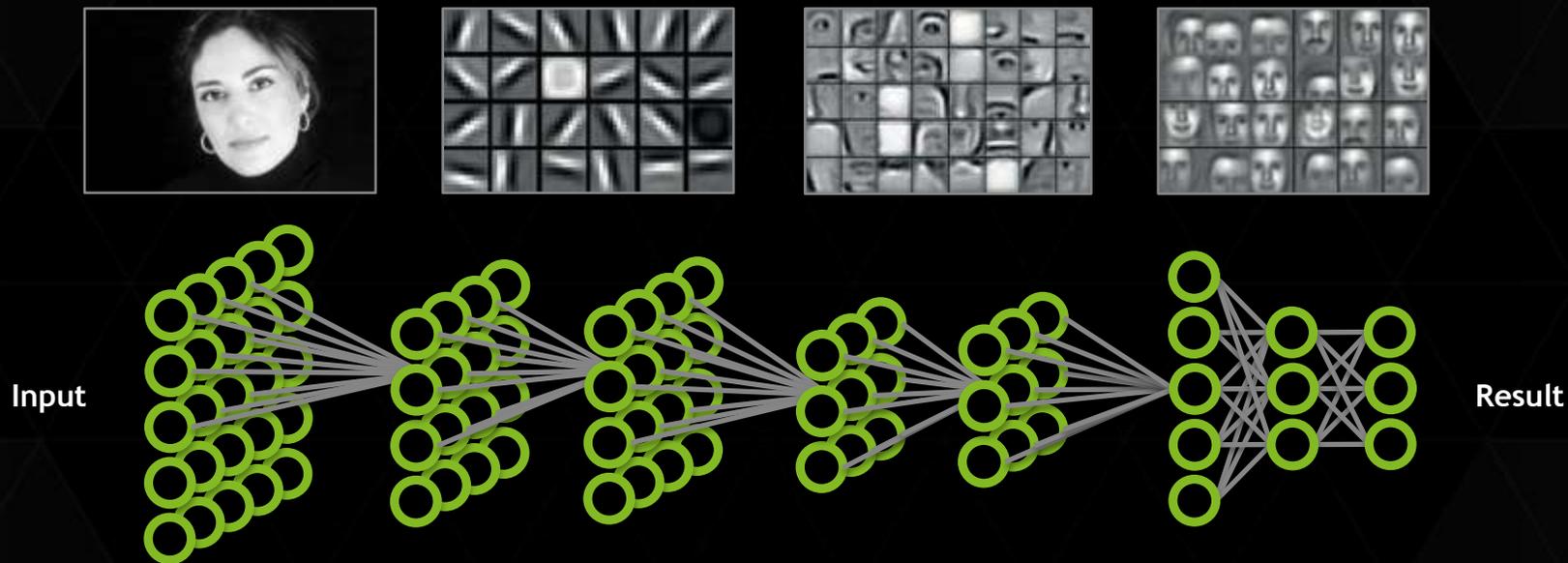


DEEP LEARNING

DEEP LEARNING FOR IMAGE ANALYTICS



MACHINE LEARNING USING DEEP NEURAL NETWORKS



Hinton et al., 2006; Bengio et al., 2007; Bengio & LeCun, 2007; Lee et al., 2008; 2009

Visual Object Recognition Using Deep Convolutional Neural Networks

Rob Fergus (New York University / Facebook) <http://on-demand-gtc.gputechconf.com/gtcnew/on-demand-gtc.php#2985>

BROAD BENEFITS OF DEEP LEARNING



Spotify

Content-based music recommendation

Summer Intern implements recommendation system



Mitosis Detection in Breast Cancer Histology Images with Deep Neural Networks

Dan C. Cirescan et al.

Research team takes first step to bring automated mitosis detection into clinical practice



Merck Molecular Activity Challenge

Winning team dominates the competition by using deep learning algorithms running on GPUs

BROAD USE OF GPUS IN DEEP LEARNING

Early Adopters



Adobe

Image Analytics
for Creative
Cloud



Speech/Image
Recognition

flickr

Image
Classification

IBM

Hadoop

NETFLIX

Recommendation

Yandex

Search Rankings

Use Cases

Image Detection

Face Recognition

Gesture Recognition

Video Search & Analytics

Speech Recognition & Translation

Recommendation Engines

Indexing & Search

Talks @ GTC

facebook



STANFORD
UNIVERSITY



DENSO

Carnegie
Mellon
University

MIT
Massachusetts
Institute of
Technology

Berkeley
UNIVERSITY OF CALIFORNIA

WHAT IS NEXT?

Deep Learning Will Be Everywhere

Pattern Analysis



Anomaly Detection



Behavior Prediction



Diagnostic Support



Sentiment Analysis

....

“Mark Zuckerberg calls it the theory of the mind. How do we model – in machines – what human users are interested in and are going to do?”

Yann Lecun, Director AI Research at Facebook

“Any product that excites you over the next five years and makes you think: ‘That is magical, how did they do that?’, is probably based on this [deep learning].”

Steve Jurvetson, Partner DFJ Venture

GPU TECHNOLOGY CONFERENCE

March 17-20, 2015 | Silicon Valley
www.gputechconf.com #GTC15



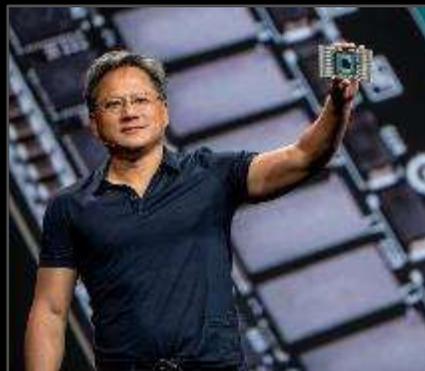
CONNECT

Connect with technology experts from NVIDIA and other leading organizations



LEARN

Gain insight and hands-on training through the hundreds of sessions and research posters



DISCOVER

See how GPU technologies are creating amazing breakthroughs in important fields such as deep learning



INNOVATE

Hear about disruptive innovations as early-stage startups present their work

2015 Theme: Deep Learning

FUELING THE DEEP LEARNING REVOLUTION

March 17 – 20, 2015 | Silicon Valley | #GTC15

REGISTER NOW

50+
Deep Learning Sessions

www.gputechconf.com

Adobe	Google
Alibaba	iFlytek, Ltd
Baidu	NUANCE
Carnegie Mellon	Stanford Univ
Facebook	UC Berkeley
Flickr / Yahoo	Univ of Toronto

Developer Labs

Caffe
Torch
Theano