



## **GPUs Today**

**Lessons from Graphics Pipeline** 

- Throughput is paramount
- Create, run, & retire lots of threads very rapidly



#### Early Electronic Graphics Hardware

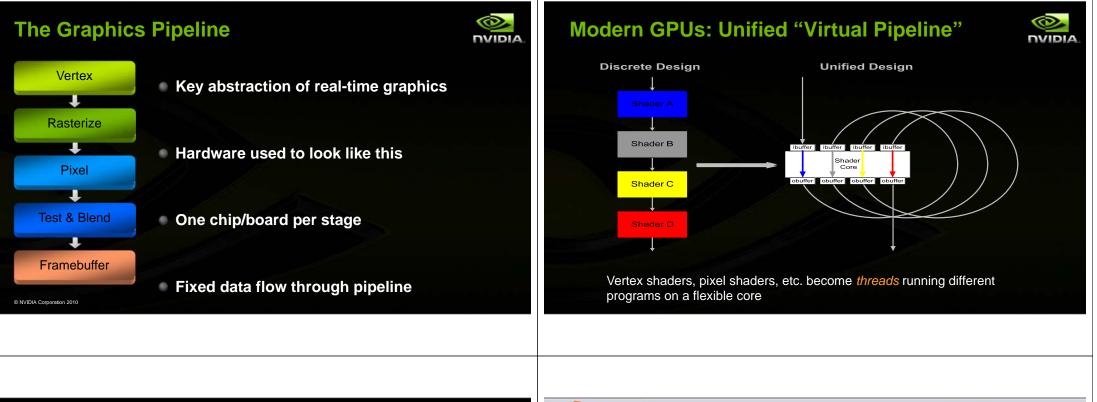




SKETCHPAD: A Man-Machine Graphical Communication System Ivan Sutherland, 1963

© NVIDIA Corporation 2010

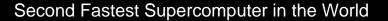
"Fermi"





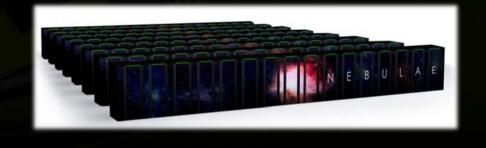


### **Dawning Nebulae**



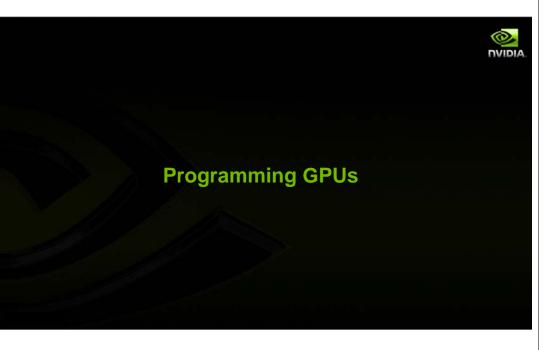
1.27 Petaflop

4640 Tesla GPUs



#### 1000+ GPU Clusters Around the World



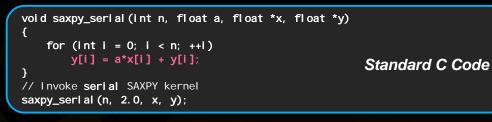




## C for CUDA : C with a few keywords



**NVIDIA** 



\_\_\_\_\_global\_\_\_ void saxpy\_parallel(int n, float a, float \*x, float \*y) {

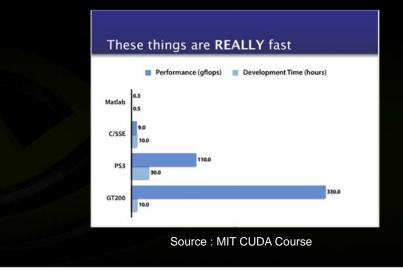
int i = blockldx.x\*blockDim.x + threadldx.x;
if (i < n) y[i] = a\*x[i] + y[i];
}</pre>

Parallel C Code

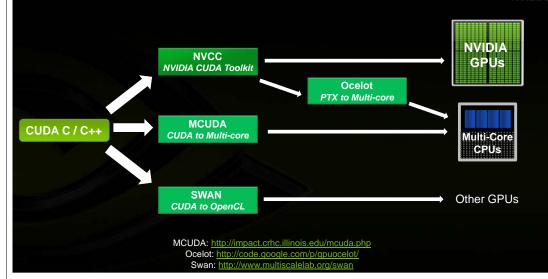
// Invoke parallel SAXPY kernel with 256 threads/block int nblocks = (n + 255) / 256; saxpy\_parallel <<<<nblocks, 256>>>(n, 2.0, x, y);



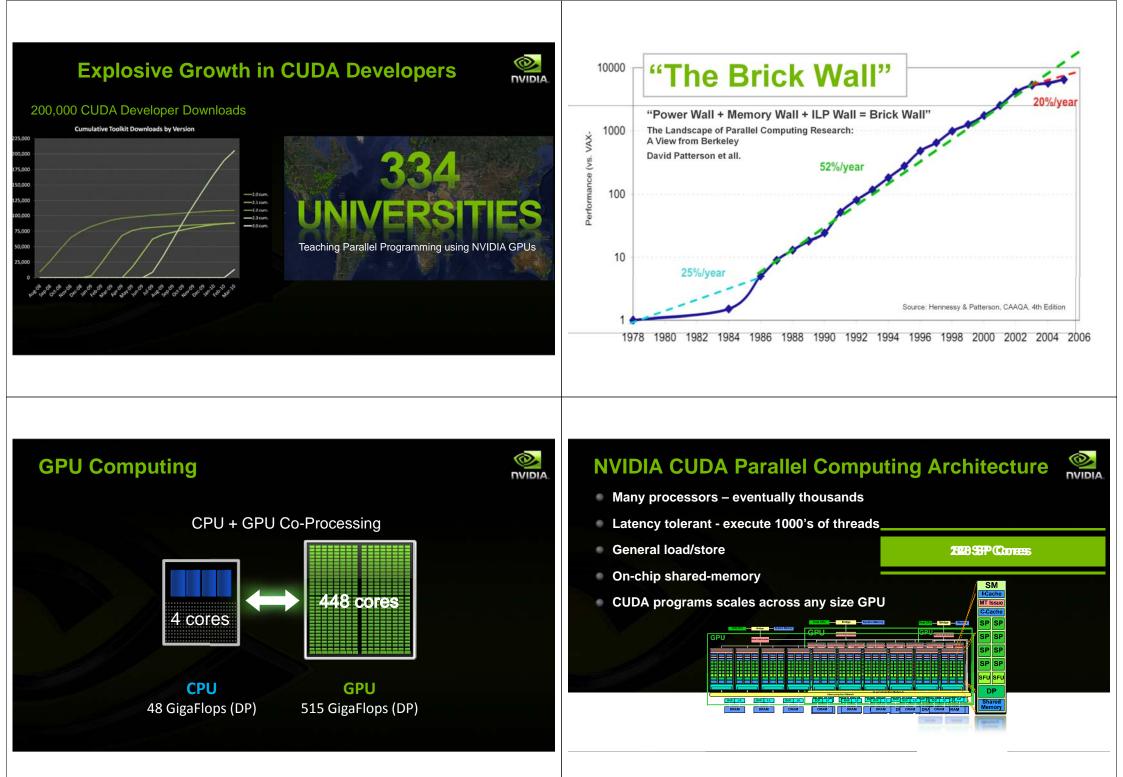
**NVIDI** 

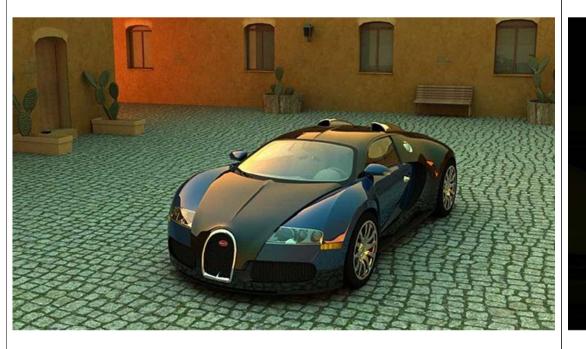


**Targeting Multiple Platforms with CUDA** 









#### **Rasterization & Ray Tracing**

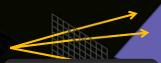
#### Rasterization

- For each triangle
  - Find the pixels it covers
  - For each pixel: compare to closest triangle so far



#### **Classical Ray Tracing**

- For each pixel
  - Find the triangles that might be closest
  - For each triangle: compute distance to pixel

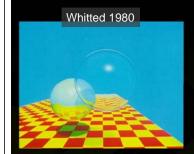


Mapped to massively parallel GPU through NVIDIA OptiX

## Why ray tracing?

- Ray tracing unifies rendering of visual phenomena
   fewer algorithms with fewer interactions between algorithms
- Easier to combine advanced visual effects robustly
  - soft shadows
  - subsurface scattering
  - indirect illumination
  - transparency
  - reflective & glossy surfaces
  - depth of field
  - ...
- But: resource intensive, challenging to make fast

## **Ray tracing regimes**



**NVIDIA** 

Mirror reflectionsPerfect refractionsHard shadows

•2-20 rays per pixel



Depth of field
Motion blur
Soft shadows
Glossy reflections
20-200 rays per pixel



**NVIDIA** 

Indirect illuminationCausticsPhysical accuracy

•200-10<sup>5</sup> rays per pixel

## **OptiX Examples**



Interactive



## **OptiX Examples**



Interactive

#### Progressive

#### **Stochastic Rasterization**

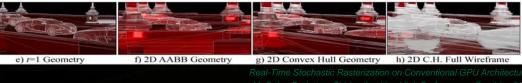


- Rasterize convex hull of time-continuous triangle
- Ray trace against TCT at each pixel



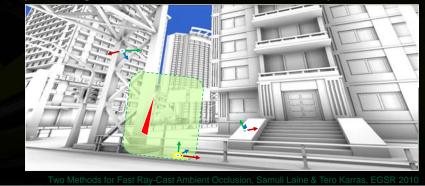






## **Ambient Occlusion**

- Darken pixels by % of hemisphere blocked by nearby triangles
- Compute triangle regions of influence to find affected pixels





#### Workloads



- Each GPU is designed to target a mix of known and speculative workloads
- The art of GPU design is choosing these workloads (and shipping on schedule!)

What workloads will drive future GPUs?

- High performance computing
- Graphics
- BIG Science
- Computational graphics

#### Histogram

- Distribution of colors in an image
- Image analysis for High Dynamic Range tone mapping





NVIDIA

Reinhard HDR tone mapping

HDR in Valve's source engine

#### **Separable Filters**

Depth of field, film bloom





lalo 3 © Bungie Studios

Crysis © Crytek Gmb

#### **Separable Filters**

Subsurface scattering via texture space diffusion





Realistic Skin Renderin Eugene d'Eon, David Luebke, Eric Enderto



#### **More Post-Processing Effects**

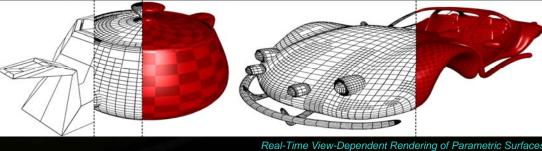




© Per Lonroth, Mattias Unger, DICE

#### **CUDA Tessellation**

- Flexible adaptive geometry generation
- Recursive subdivision



Eisenacher, Meyer, Loop 2009



#### **Real-time fluid effects**

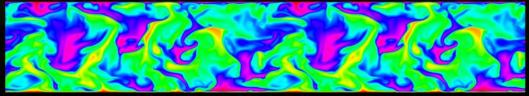
© NVIDIA Corporation 2010

Complex fluid-drive motion is all around

 Car exhaust, dust storms, rolling mist, steam, smoke, fire, contrails, bubbles in water, ...

Goal: Add this level of realism to games

Problem: Turbulent motion is computationally intensive!





NVIDIA

#### Solution: GPUs are computational monsters!

Calculate near-field fluid on grid Fluid velocities drive particle • motions

- 1. Calculate Fluid Velocities on Regular Grid 2<sup>nd</sup>-Order Accurate CUDA Multigrid Solver
- 2. Interpolate Fluid Velocities Conto Particles 3D Interpolation in CUDA
- 3. Advance Particles CUDA Particle System
- 4. Render Particles CUDA - OpenGL Interop



#### **APEX Turbulence**

#### Interactive CFD Solution + Volume Rendering







NVIDIA

#### Making Science Better, not just Faster

Already Available					NVIDIA.				
				Q	2	C	Q4		
Tools &	CULA CUDA C/C++, LAPACK Library PGI Fortran		Nsight Visual Studio IDE	Studio IDE Debugger			PGI Accelerator Enhancements	s	
Libraries	Thrust: C++ Template Lib	Jacket: MATLAB Plugin	NPP Performance Primitives (NPP)	Platform Cluster Management	Bright Cluster Management	Mathematica	CAPS HMPP Enhancements	Mathworks MATLAB	
Oil & Gas	Seismic Analysis: ffA, HeadWave	Seismic Analysis: Geostar	Seismic City, Acceleware			Seismic Interpretation	Reservoir Simulation 1	Reservoir Simulation 2	
Bio- Chemistry	AMBER, GROM HOOMD, LAMM	IACS, GROMOS, PS, NAMD, VMD	BigDFT, ABINIT, TeraChem		Quantum Chem Code 1	Other Popular MD code	Quantum Chem Code 2		
Bio- Informatics	Hex Protein Docking	CUDA-BLASTP, MUMmerGPU, M		Protein Docking		Short-read seq analysis			
Video & Rendering	Fraunhofer JPEG2K	OptiX Ray Tracing Engine	mental ray with iray	Main Concept Video Encoder	Elemental Video Live	3D CAD SW with iray			
Finance	NAG: RNGs	NumeriX: CounterParty Risk	Scicomp SciFinance	Risk Analysis 1	Risk Analysis 2	Credit Risk Analysis ISV	Trading Platform ISV		
CAE	AutoDesk Moldflow	OpenCurrent: CFD/PDE Library	Moldex3D	Acusim AcuSolve CFD	Structural Mechanics ISV	MSC MARC	MSC Nastran	Several CAE ISVs	
EDA	Electro-magnetics: Agilent, CST, Remcom, SPEAG S		Agilent ADS Spice Simulator			Verilog Simulator	Lithography Products	SPICE Simulator	
			d the						
					Delessed				

**Increasing Number of CUDA Applications** 



Product

	An Exc						
146X	36X	18X	50X	100X	6:0	Q: ()	\$; (E) =
Medical Imaging U of Utah	Molecular Dynamics U of Illinois, Urbana	Video Transcoding Elemental Tech	Matlab Computing AccelerEyes	Astrophysics RIKEN			(
	5	50x – 150				Image	
149X	47X	20X	130X	30X		• A	odium ny sig uch le ery lar
Financial simulation	Linear Algebra	3D Ultrasound	Quantum Chemistry	Gene Sequencing		R	equire

U of Illinois, Urbana

#### iting Revolution - Sodium Map of the Brai

				۲			٢		0		0	1
	X	$\langle \mathbf{\hat{x}} \rangle$								, act		÷.
3:	Ś	3	4		1 1 2	$\mathcal{P}_{\mathbf{r}}$	1990 A					

Courtesy of Keith Thulborn and Ian Atkinson, Center for MR Research, University of Illinois at Chicago

#### sodium in the brain

- m is one of the most regulated substance in human tissues
- ignificant shift in sodium concentration signals cell death
- less abundant than water in human tissues, about 1/2000
- arge number of samples are needed for good SNR
- res high-quality reconstruction, currently considered impractical

Thanks: Wen-mei Hwu

U of Maryland

An Exciting Revolution - Sodium Map of the Brai

Techniscan

Courtesy of Keith Thulborn and Ian Atkinson, Center for MR Research, University of Illinois at Chicago

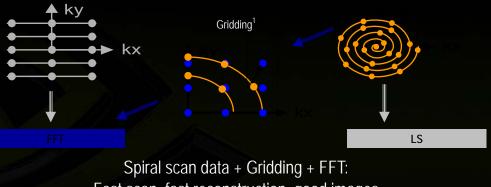
- Enables study of brain-cell viability before anatomic changes occur in stroke and cancer treatment.
  - Drastic improvement of timeliness of treatment decision
  - Minutes for stroke and days for oncology. ۰

Universidad Jaime

Oxford

# **Reconstructing MR Images**

Spiral Scan Data



Fast scan, fast reconstruction, good images Can become realtime with about 10X speedup.

<sup>1</sup> Based on Fig 1 of Lustig et al, Fast Spiral Fourier Transform for Iterative MR Image Reconstruction, IEEE Int'l Symp. on Biomedical Imaging, 2004

## 

Spiral scan data + LS Superior images at expense of significantly more computation; several hundred times slower than gridding. Traditionally considered impractical!

kx

#### **Summary of Spiral Scan LS Results**

						DVIDIA.
	Q	Q				
Reconstruction	Run Time (m)	GFLOP	Run Time (m)	GFLOP	Linear Solver (m)	Recon. Time (m)
Gridding + FFT (CPU, DP)	N/A	N/A	N/A	N/A	N/A	0.39
LS (CPU, DP)	4009.0	0.3	518.0	0.4	1.59	519.59
LS (CPU, SP)	2678.7	0.5	342.3	0.7	1.61	343.91
LS (G80, Naïve)	260.2	5.1	41.0	5.4	1.65	42.65
LS (G80, CMem)	72.0	18.6	9.8	22.8	1.57	11.37
LS (G80, CMem, SFU)	13.6	98.2	2.4	92.2	1.60	4.00
LS (G80, CMem, SFU, Exp/layout)	7.5	178.9	1.5	) 144.5	1.69	3.19
NVIDIA Confidential	<ul> <li>357X, in mach up time.</li> </ul>	nine set	228X -			

High-Throughput Computing = Futuristic Biology

- in-silico screening of drugs
- mastering diseases

FFT

personalized medicine

## **In-silico Drug Screening**

- Weed out inactive compounds
- Rank "drug candidates" for given targets
- Example:

O,

NVIDIA

- CERN grid 300,000 potential drugs against avian flu screened
- 2000 computers, 4 weeks!
- 4 years cpu-time

## protein aggregation

O.

a process critical in

- some degenerative diseases (e.g., Parkinson's): aggregates abnormal
- drug production: aggregates undesirable

time scale of the process:

- → in vitro: up to days!
- impossible for molecular dynamics

Thanks: Lorena Barba

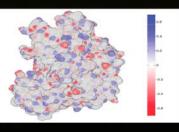
#### **Electrostatic Interactions Play a Crucial Role**



Classical molecular dynamics:

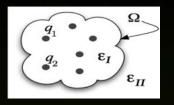
very detailed ... but too expensive at large scale!

- Alternative: continuum model of surrounding water
  - don't care what the H<sub>2</sub>0 molecules do
  - model as a continuum dielectric
  - leads to a boundary integral equation (BIE) problem
- Fast algorithm, well-suited for GPU:
  - fast multipole method, solves BIE in O(N) ops



#### As in Many Computation-hungry Applications

- Three-step approach:
  - 1. Restructure the mathematical formulation
- 2. Innovate at the algorithm level
- 3. Tune core software for hardware architecture



## Vision—predictive biology, faster, cheaper, accurate



- Drug screening:
  - Few weeks, on a (say) 32-node GPU cluster >> safe drug to market
- Protein aggregation:
  - Get physics right + 100x larger simulation >> understand Parkinson's
- Analogy:
  - circuit design : it is all done digitally and verified; the circuit works!
  - If it didn't work, it would be too costly for many consumer electronics

#### **Conclusion: Three Options**

- Accelerate Legacy Algorithms and Applications
  - Use libraries, recode existing apps
     => good work for domain scientists (minimal CS required)
- Rewrite / Create new Approaches
  - Opportunity for clever algorithmic thinking
    - => good work for computer scientists (minimal domain knowledge required)

#### Rethink Numerical Methods & Algorithms

- Potential for biggest performance advantage
  - => Interdisciplinary: requires CS and domain insight
  - => Exciting time to be a computational scientist



#### Key GPU Workloads

- Computational graphics (don't forget DirectXn)
- Scientific and numeric computing
- Image processing video & images
- Computer vision / Computational Photography
- Speech & natural language
- Data mining & machine learning



#### **Final Thoughts – Education**

We should teach parallel computing in CS 1 or CS 2
Computers don't get faster, just wider
now
Manycore is the force of computing... and graphics

Insertion Sort Heap Sort Merge Sort
Which goes faster on large data?
ALL Students need to understand this! Early!

**NVIDIA**