



SIGGRAPH2007

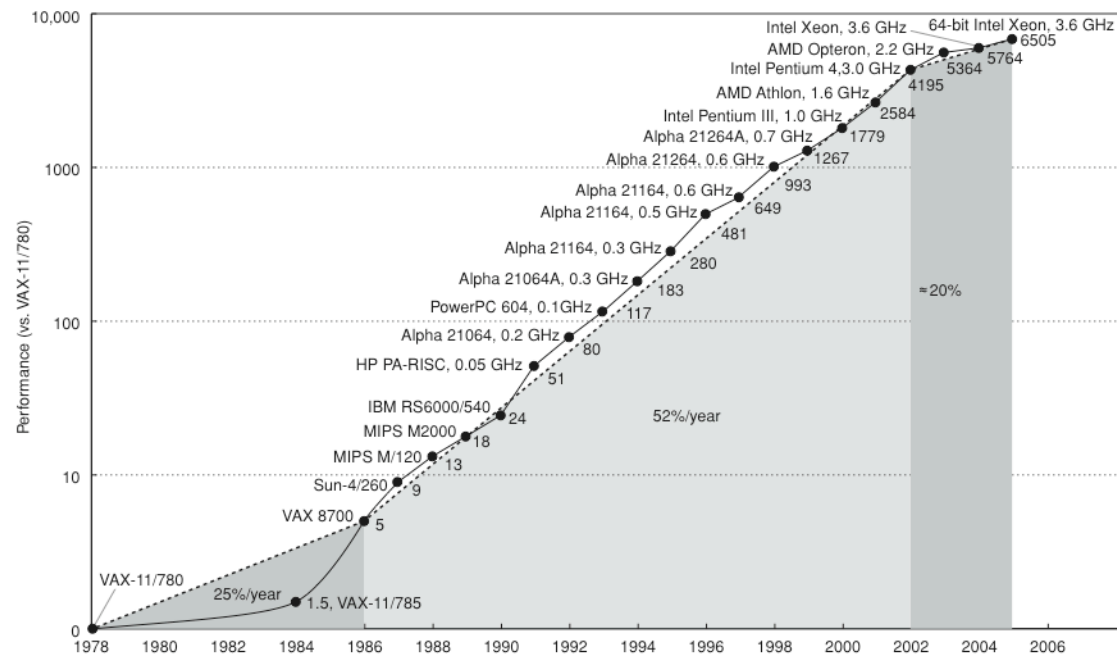
GPU Architecture Overview

John Owens

UC Davis



The Right-Hand Turn



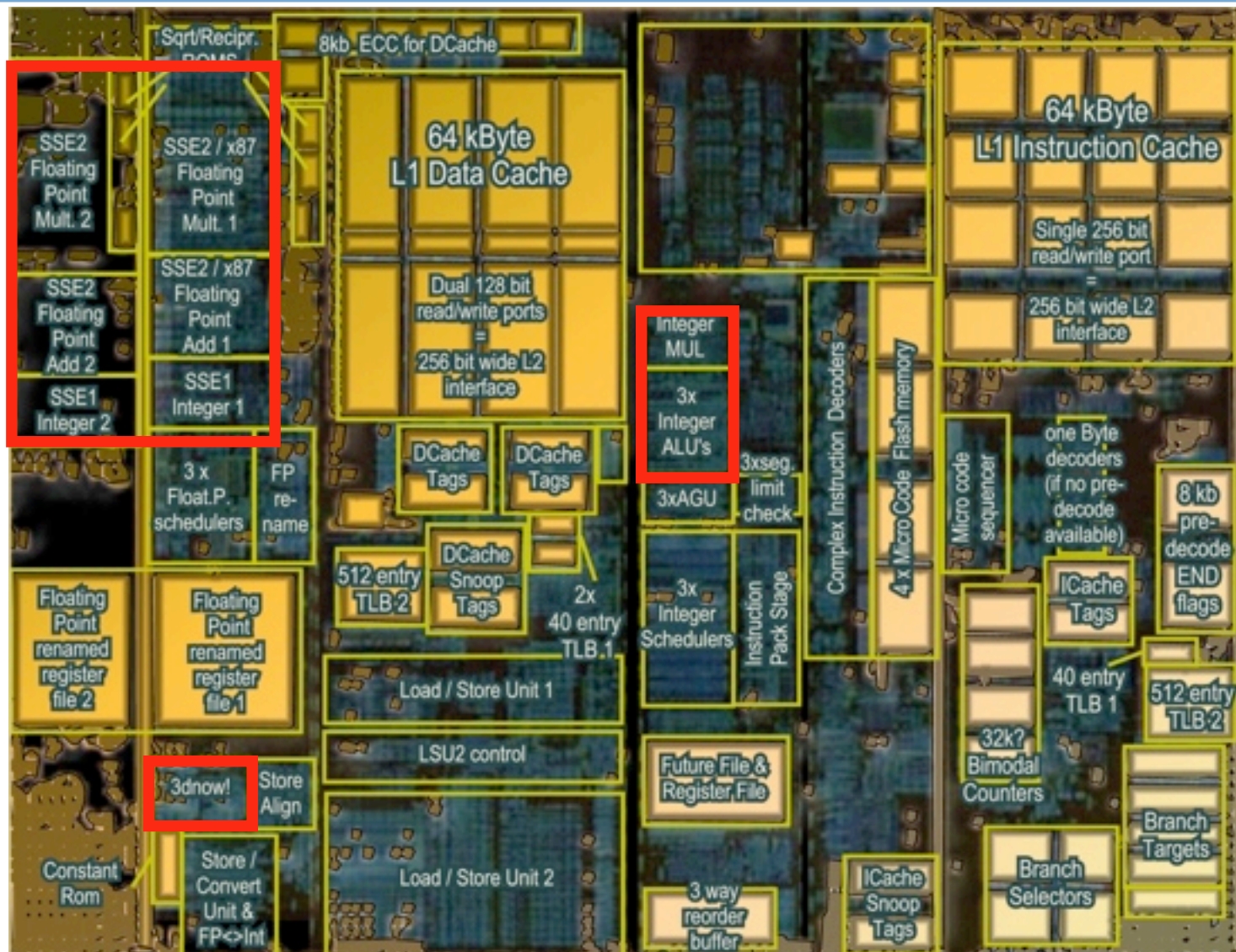
Why? [Architecture Reasons]

- ILP increasingly difficult to extract from instruction stream
- **Control hardware dominates μ processors**
 - Complex, difficult to build and verify
 - Takes substantial fraction of die
 - Scales poorly
 - Pay for max throughput, sustain average throughput
 - Quadratic dependency checking
 - Control hardware doesn't do any math!
 - Intel Core Duo: 48 GFLOPS, ~10 GB/s
 - NVIDIA G80: 330 GFLOPS, 80+ GB/s

AMD “Deerhound” (K8L)



SIGGRAPH2007



Why? [Technology Reasons]

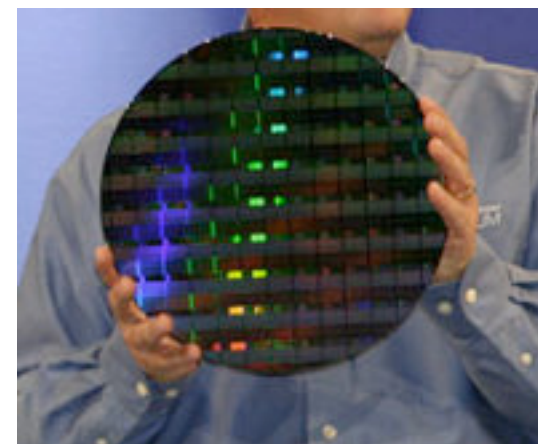
- **Industry moving from “instructions per second” to “instructions per watt”**
 - “Power wall” now all-important
 - Traditional μ proc techniques are not power-efficient
- **We can continue to put more transistors on a chip ...**
 - ... but we can’t scale their voltage like we used to ...
 - ... and we can’t clock them as fast ...



SIGGRAPH2007

Go Parallel

- **Time of architectural innovation**
 - GPUs let us explore using hundreds of processors now, not 10 years from now
- **Major CPU vendors supporting multicore**
- **Interest in general-purpose programmability on GPUs**
- **Universities must teach thinking in parallel**



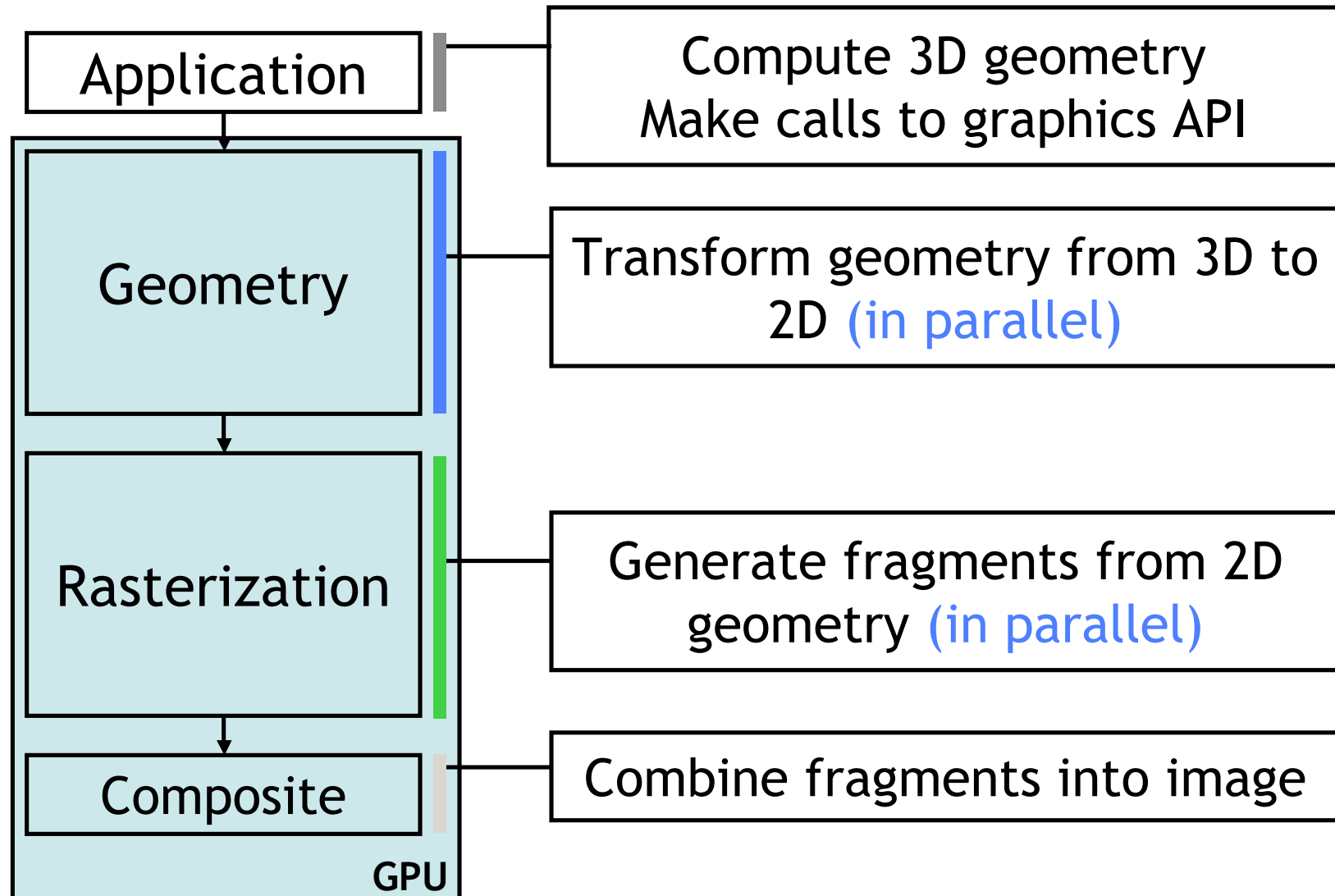
What's Different about the GPU?

- **The future of the desktop is parallel**
 - We just don't know what kind of parallel
- **GPUs and multicore are different**
 - Multicore: Coarse, heavyweight threads, better performance per thread
 - GPUs: Fine, lightweight threads, single-thread performance is poor
- **A case for the GPU**
 - Interaction with the world is visual
 - GPUs have a well-established programming model
 - Market for GPUs is 500M+ total/year

The Rendering Pipeline



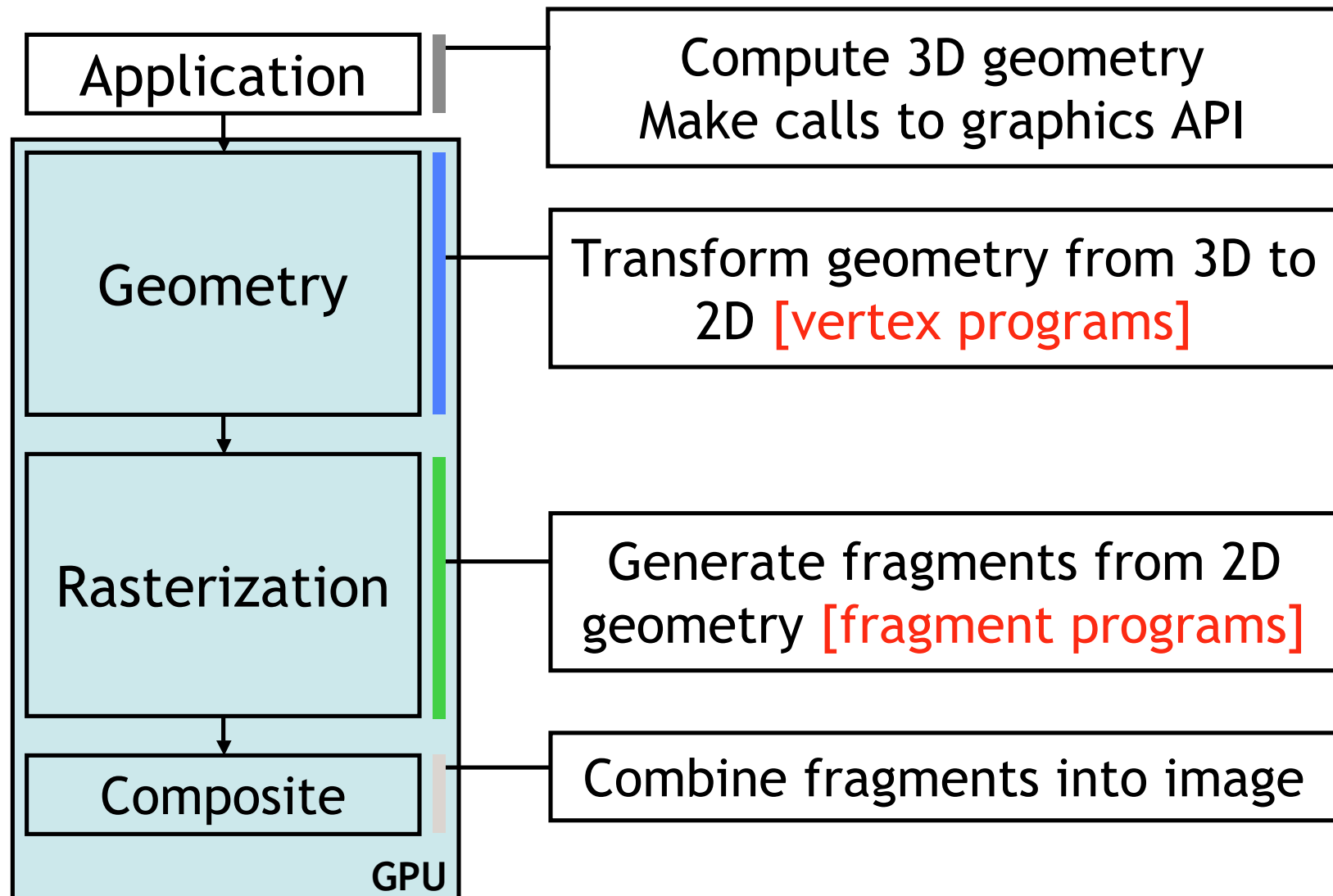
SIGGRAPH2007



The Programmable Pipeline



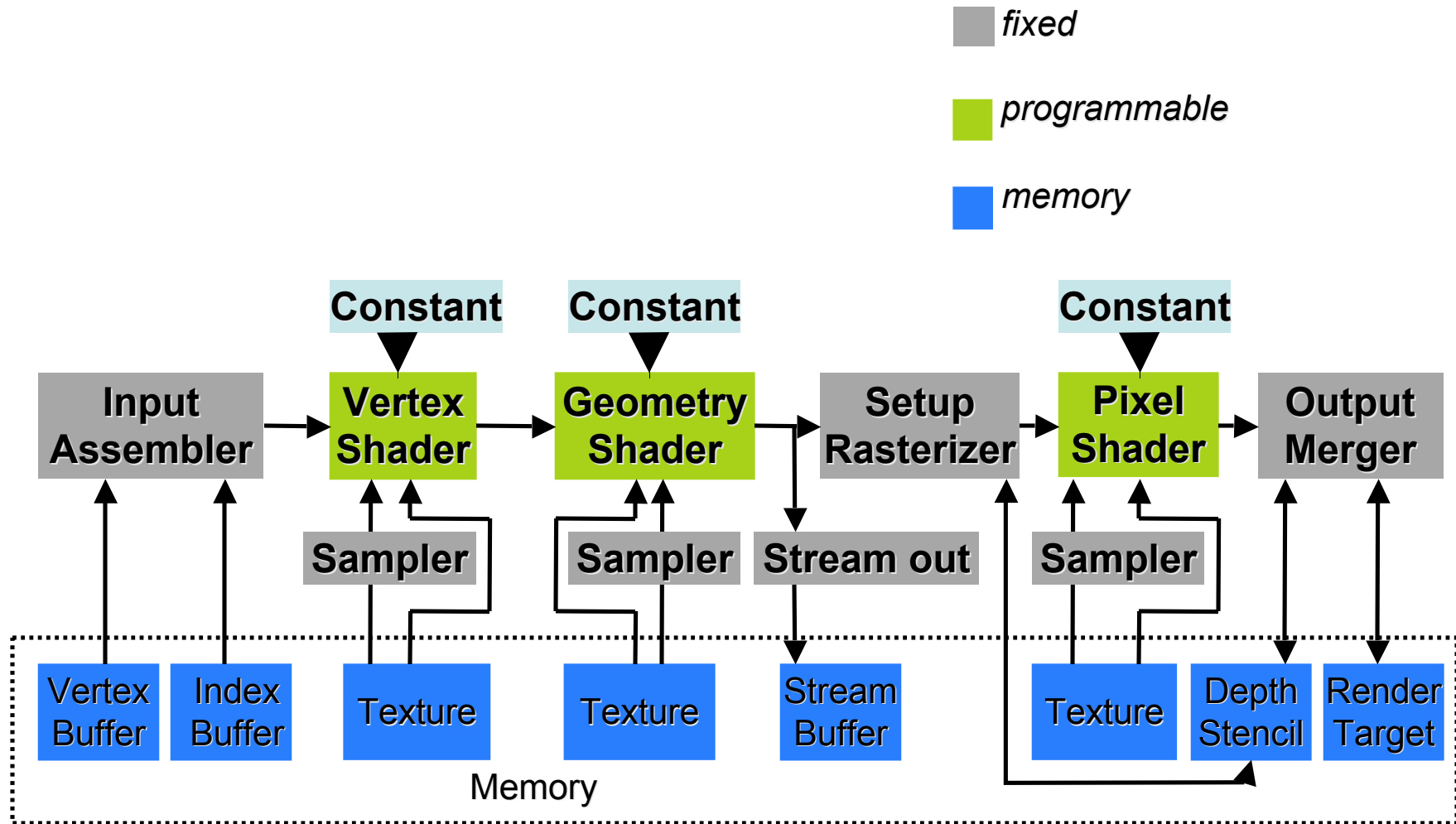
SIGGRAPH2007



DirectX 10 Pipeline



SIGGRAPH2007



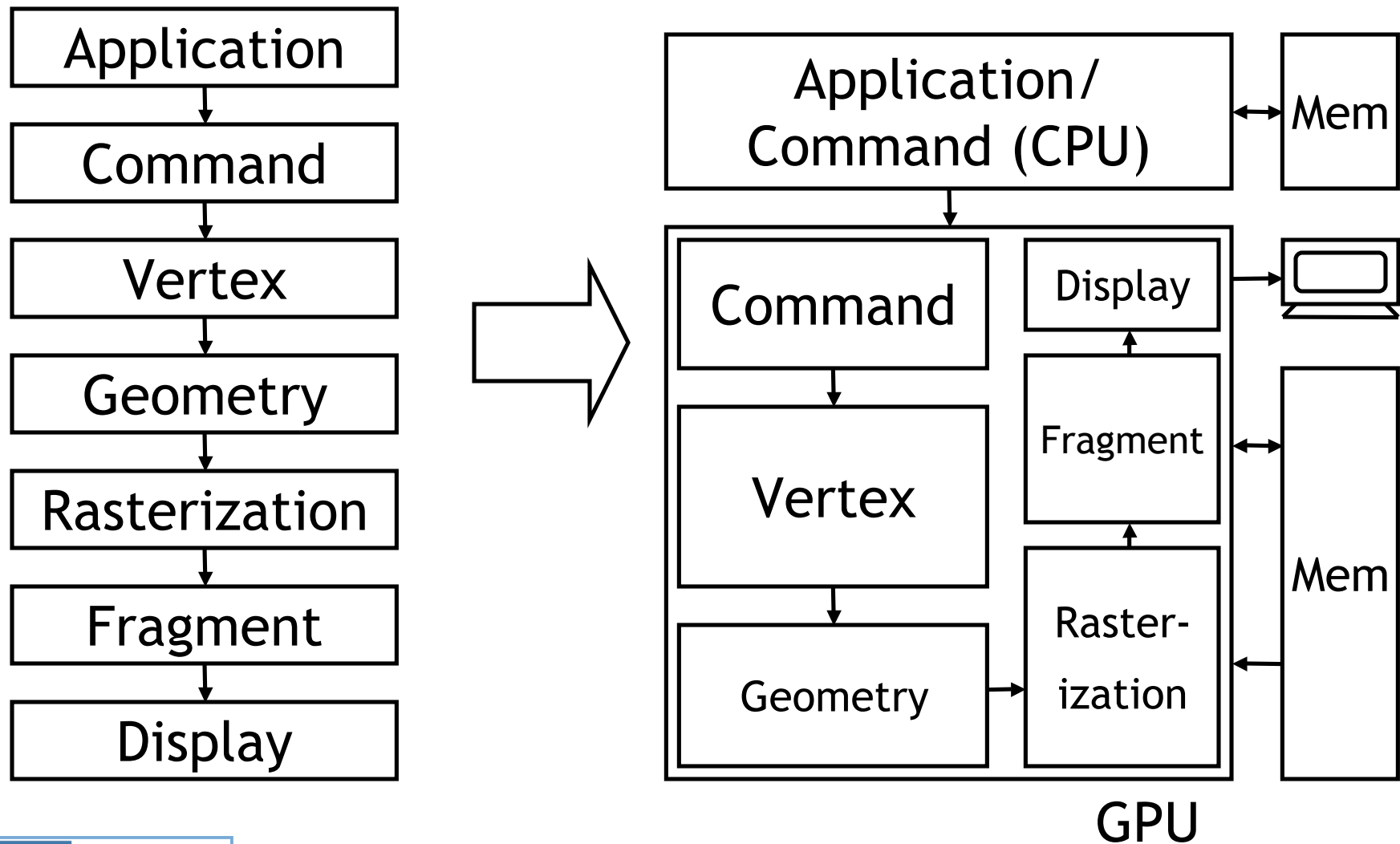
Characteristics of Graphics

- **Large computational requirements**
- **Massive parallelism**
 - Graphics pipeline designed for independent operations
- **Long latencies tolerable**
- **Deep, feed-forward pipelines**
- **Hacks are OK—can tolerate lack of accuracy**
- **GPUs are good at parallel, arithmetically intense, streaming-memory problems**

Graphics Hardware—Task Parallel



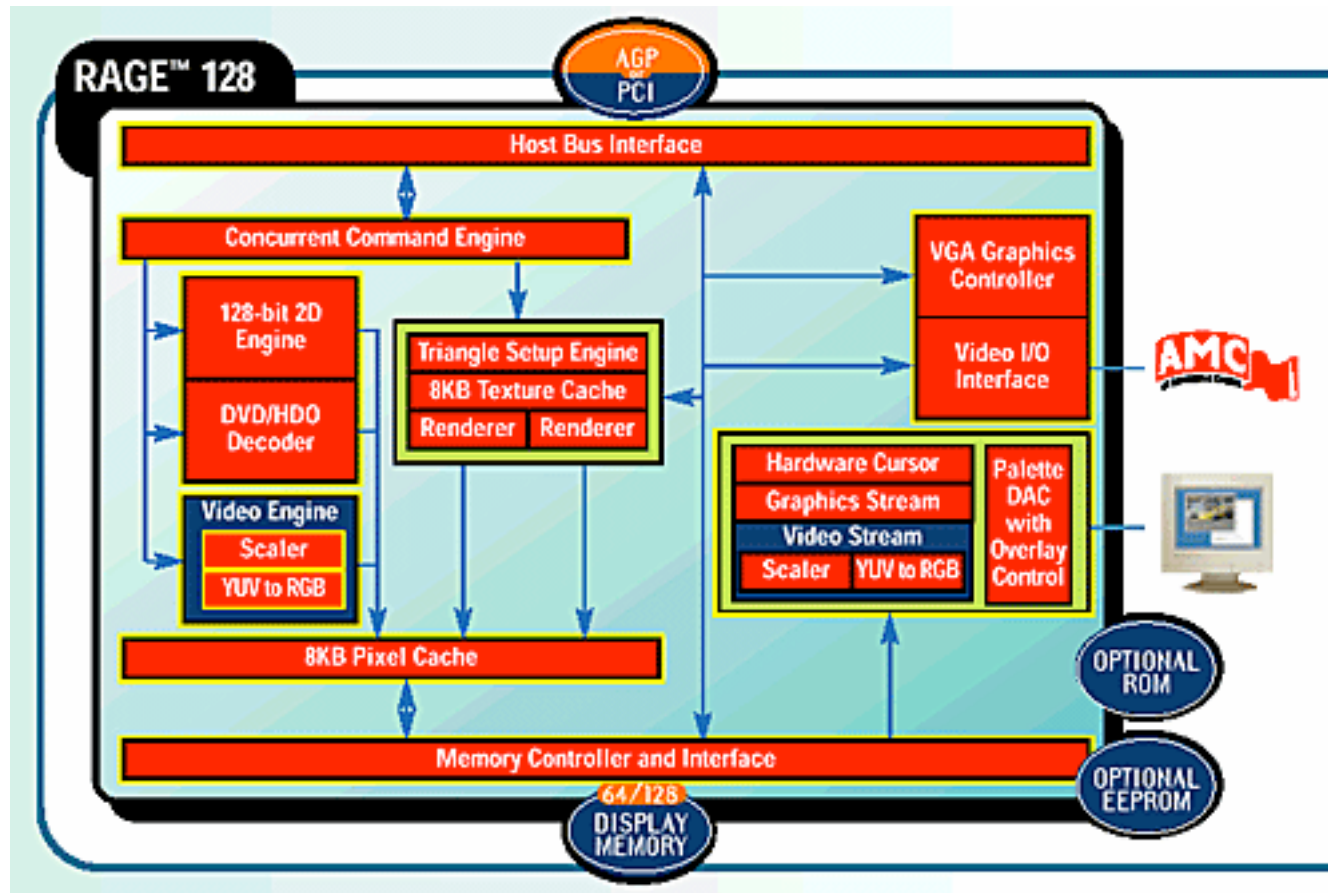
SIGGRAPH2007



Rage 128



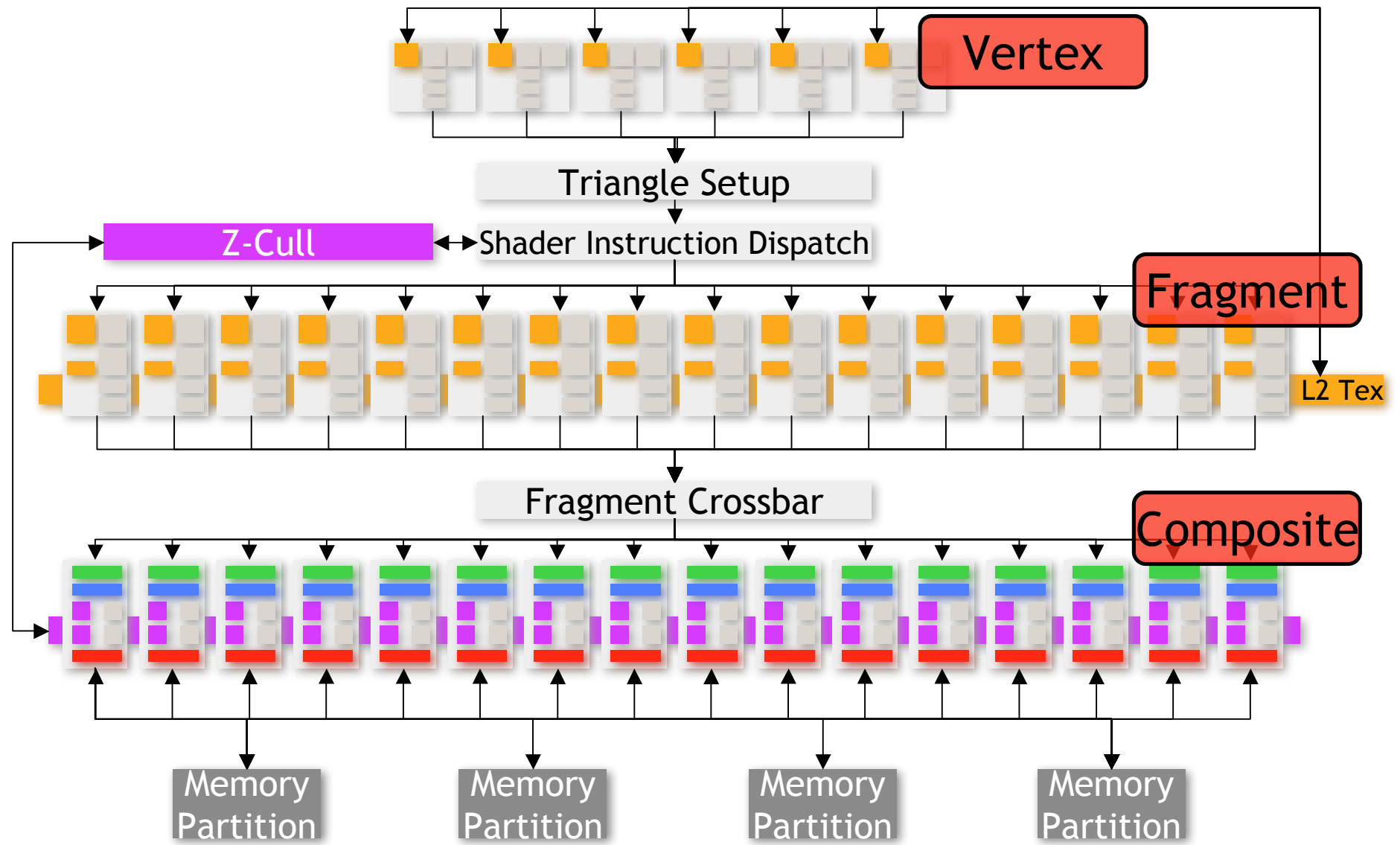
SIGGRAPH2007



NVIDIA GeForce 6800 3D Pipeline



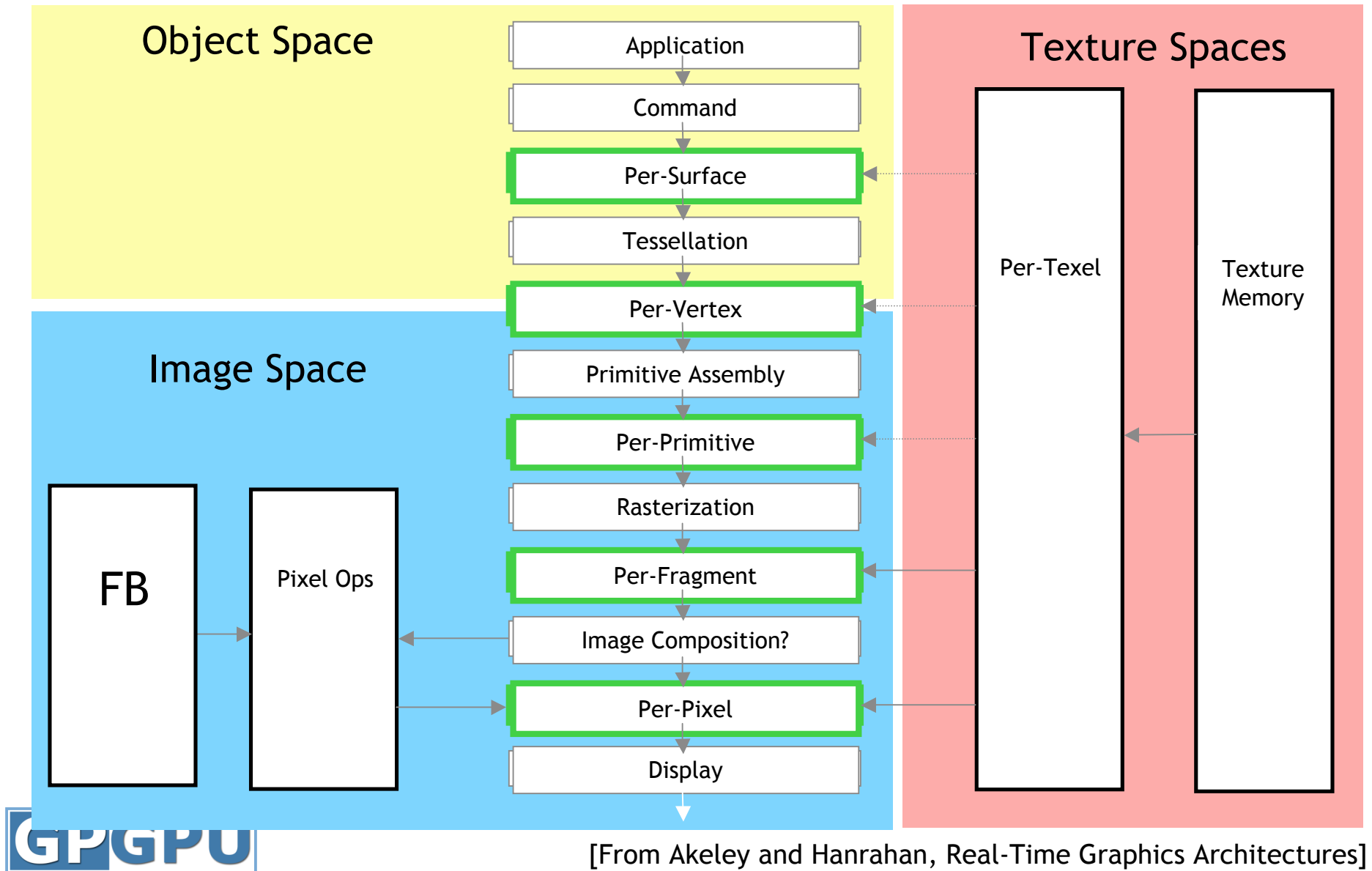
SIGGRAPH2007



Programmable Pipeline



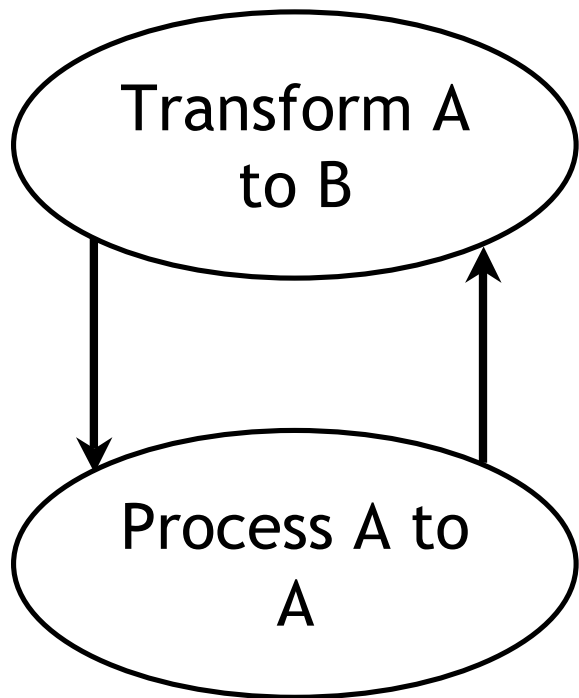
SIGGRAPH2007



Generalizing the Pipeline



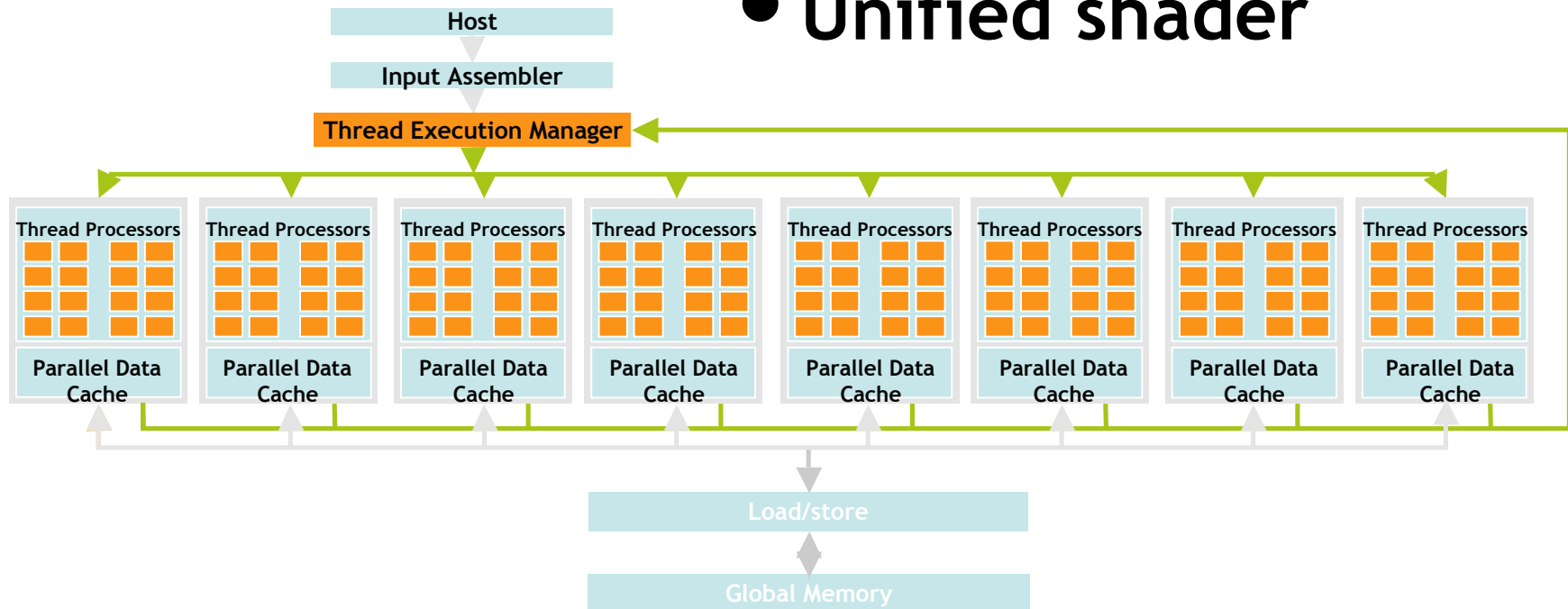
SIGGRAPH2007



- **Transform A to B**
 - Ex: Rasterization (triangles to fragments)
 - Historically fixed function
- **Process A to A**
 - Ex: Fragment program
 - Recently programmable, and becoming more so

GeForce 8800 GPU

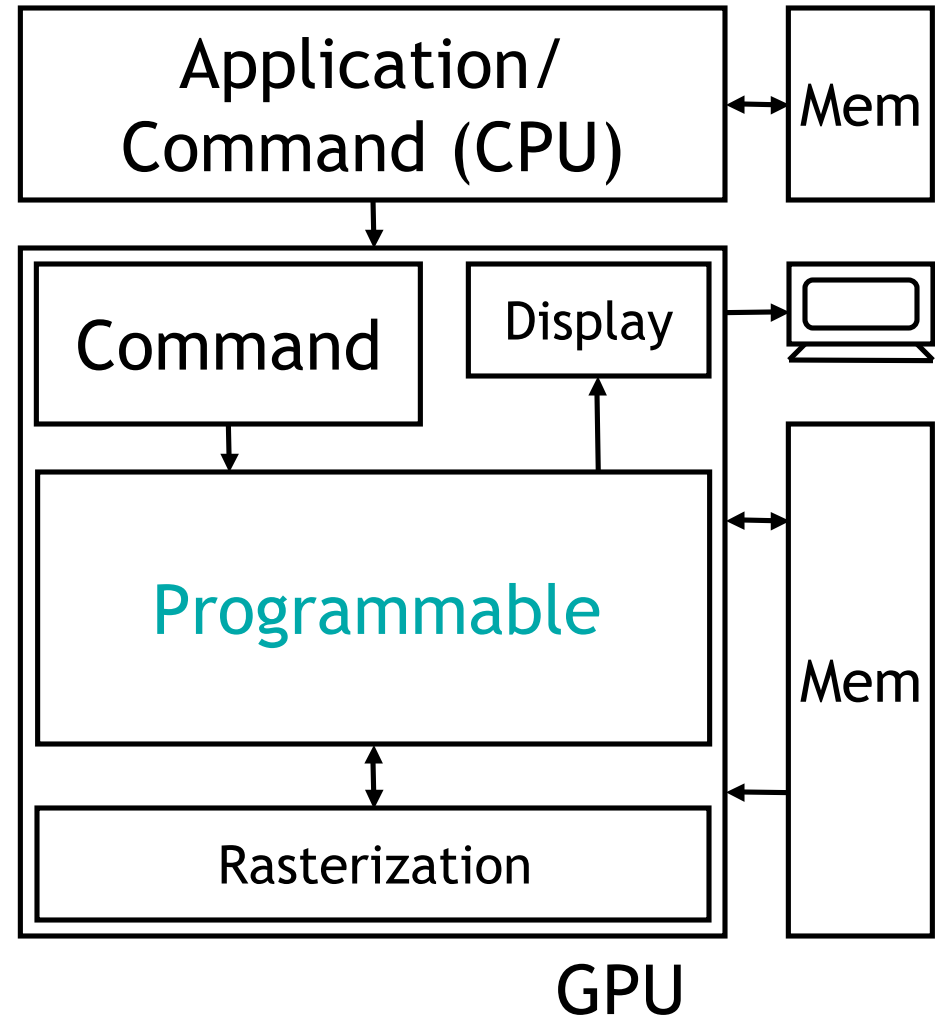
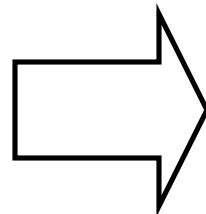
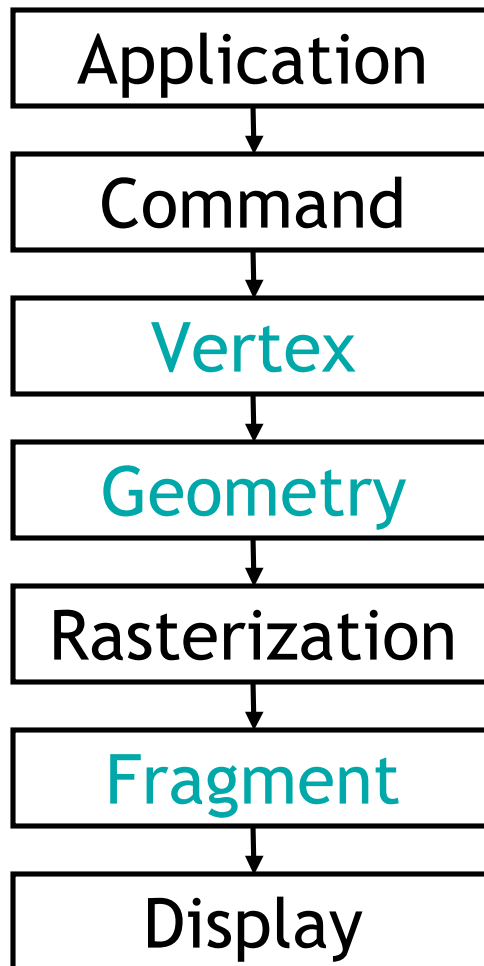
- Built around programmable units
- Unified shader



Unified Shaders



SIGGRAPH2007



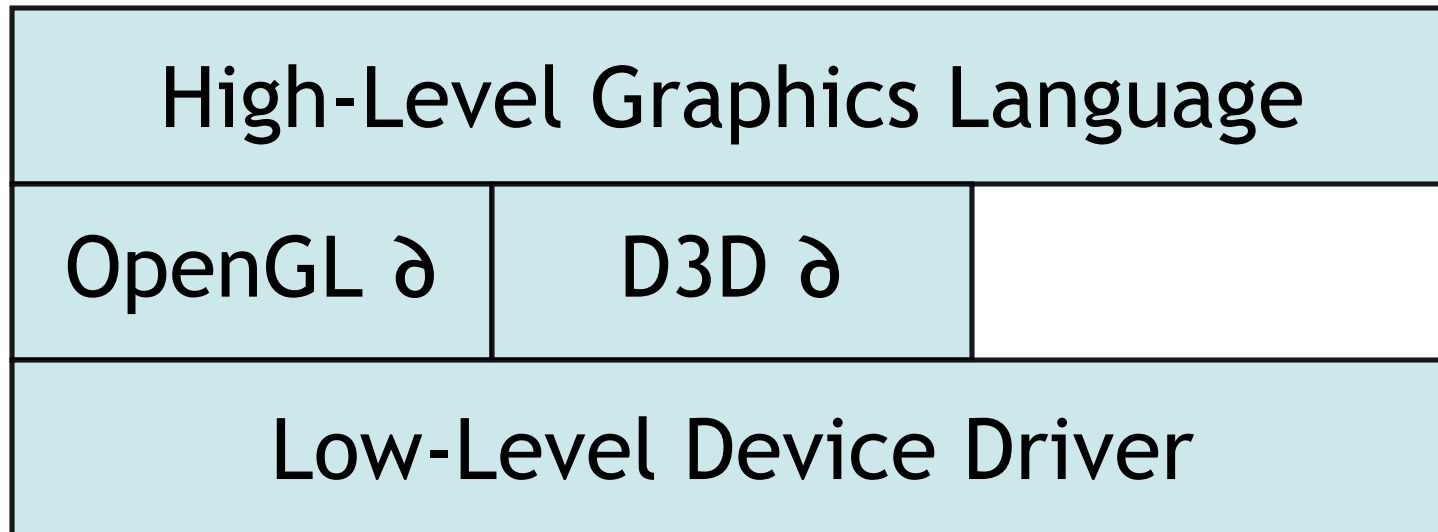
Towards Programmable Graphics

- **Fixed function**
 - Configurable, but not programmable
- **Programmable shading**
 - Shader-centric
 - Programmable shaders, but fixed pipeline
- **Programmable graphics**
 - Customize the pipeline
 - Neoptica asserts the major obstacle is programming models and tools

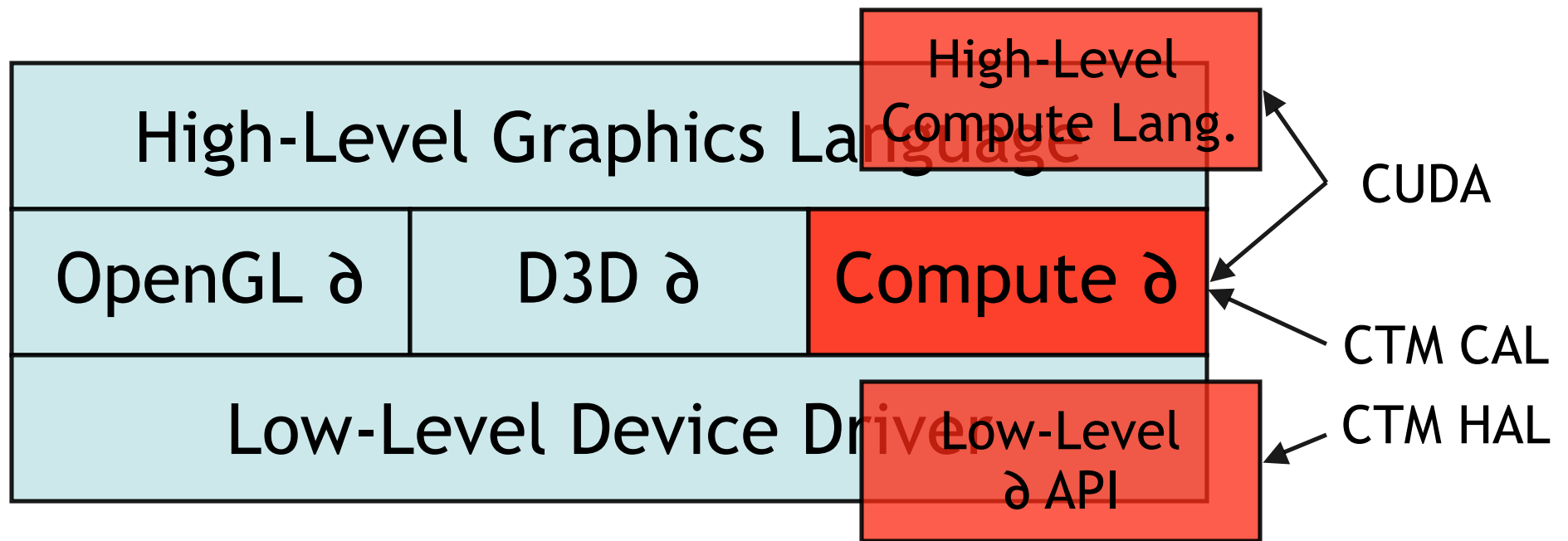
<http://www.neoptica.com/NeopticaWhitepaper.pdf>

http://www.graphicshardware.org/previous/www_2006/presentations/pharr-keynote-gh06.pdf

Yesterday's Vendor Support



Today's New Vendor Support



Architecture Summary



SIGGRAPH2007

- **GPU is a massively parallel architecture**
 - Many problems map well to GPU-style computing
 - GPUs have large amount of arithmetic capability
 - Increasing amount of programmability in the pipeline
- **New features map well to GPGPU**
 - Unified shaders
 - Direct access to compute units in new APIs
- **Challenge:**
 - How do we make the best use of GPU hardware?
 - Techniques, programming models, languages, evaluation tools ...