

Architectural Drivers

Programmability
 Parallelism
 Memory bandwidth



1

Recent History: GeForce 1&2

- First integrated geometry engine & 4 pixels/clk
- Fixed-function transform, lighting, and pixel pipelines
- 25M transistors : 0.18um/6LM : 250MHz
- 25M polygons/sec : 1G pixels/sec





Rendering in Transition

Pre-2001: pixel "painting"

- Image complexity and richness from LOTS of pixels
- Each pixel derived from 1-2 textures & blending
- Detail added by transparency and layers
- Post-2001 fork in the road:
 - Paint more simple pixels, faster embedded DRAM OR
 - Use Programmable Shading to render "better" pixels - but, must reduce depth complexity



GPUs vs. CPUs

More independent calculations

 Enables wide and deep parallelism
 API churn

 shorter development cycles -> ASIC

 Blend of general- and special- purpose compute resources
 Both transistor-bound for the foreseeable

Both transistor-bound for the foreseeable future



Special-Purpose Hardware

Most efficient implementations of

- Cube environment map
- Shadow calculations
- Anisotropic filtering
- Clipping
- Rasterization
- Log, exp, dot-product

More programmability won't change this



Managing DRAM Bandwidth

- Very large working sets
- Affordable caches cannot support long-term reuse
- Target is effective streaming with local reuse
- Supercomputer techniques apply
 - Latency hiding
 - Vector operations



Fit the Machine to DRAM Characteristics

Page locality

- ORAMs pages are 1-D
- Graphics accesses are 1-D, 2-D, 3-D in memory
- Page misses and read/write turns getting more expensive
- Granularity



Unified Memory

Traffic comprises: Commands Vertex data Texture samples Depth values Relative amounts vary widely Powerful programming model

primitives,state
{x,y,z,w}



Eliminate Redundant Traffic

We cache data Textures, vertices We cache work Pixel fragments Designed for 80 - 90% hit rate (not 99.9%) Leverage coherence Engines traverse locally ex: rasterization order



Amplify Peak Bandwidth

Lossless compression
 Lossy compression

 Conservative (undetectable)
 Else application must be given the choice
 Small-grained random access favors fixed compression atoms



Reduce Dead Cycles

Queuing
 Amortize read/write turns over longer runs
 Multi-bank DRAM scheduling



Embedded DRAM

Tempting to embed megabytes of DRAM.
But ..

- Cannot fit the whole problem
- Costs are huge

Will be "just around the corner" for a long time



A Tour of the GeForce4







Vertex Program Examples



NVIDI

- Deformation
- Warping

Lens Effects

Procedural Animation



- Range-based Fog
- Elevation-based Fog
- Animation
 - Morphing
 - Interpolation







Occlusion Culling & Programmable Shading

Occlusion Culling reduces Depth Complexity

- Calculate Z and determine visible pixels
- Eliminate invisible pixels

Programmable Shading enables richer visual quality

- Accurately model: reflections, shadows, materials
- More textures/pixel
- More calculations/pixel consumes many cycles
- Programmable Shading impractical without Occlusion Culling





Pixel Shading / Texturing

- A pixel shader converts texture coordinates into a color using a shader program.
 - Floating point math
 - Texture lookups
 - Results of previous pixel shaders
- 4 stages, 1 texture address op per stage
 - Compressed, mipmapped 3-D textures
 - True reflective bump mapping
 - True dependent textures (lookup tables)
 - Full 3×3 transform with cubemap or 3-D texture lookup
 - 16-bit-per-component normal maps



Pixel Shader

- Input: values interpolated across triangle
- IEEE floating point operations
- Lookup functions using textures
 - Large, multi-dimensional tables
 - Filtered
- Outputs an ARGB value that register combiners can read





- Alternatively, each stage can evaluate dot products instead of multiplies
- Can conditionally select A*B or C*D

Pixel Shading effects

- s/IDI/
- **Multi-texturing** 0
- Dot products for per pixel lighting calculations
- Reflections 0
- Shadowing
- **Custom effects** 0
- **Pixel math** 0











24







Multisample Antialiasing

Transparent to the application
2 & 4 subsamples per pixel
2, 4, 5, and 9-tap reconstruction filters



nView Display Technology

Flexible display combinations





Intuitive user interface

- Easy set-up
- Application Management
 Window Effects
- Multi Desktop Support
- Window Management
- Custom User Profiles



Framebuffer Controller 128-pin DDR Schedules requests from all engines Transparent compress/decompress Maps from pixel-linear address to page & partition tiling Flexible in:

- Width
- Depth
- Frequency
- Banks





Statistics

- 136M vertices per second
 60M triangles per second
 4.8G samples/sec
 1.2T ops/sec
 83.2 GB/sec clear BW
 63M transistors
 TSMC 0.15u
- 300 MHz pipeline / 325 MHz memory clk





Conclusions



Questions



32

Related reading

- Stanford special topics course *Real-Time Graphics Architectures* with Kurt Akeley & Pat Hanrahan
 - Reading list:

http://graphics.stanford.edu/courses/cs448a-01-fall/readings.html

