

Mattan Erez


 The University of Texas at Austin

Outline

- 3D graphics pipeline
- Programmable graphics pipeline
- General GPU architecture
 - And high-level examples
- Some slides courtesy David Kirk (NVIDIA) and Wen-Mei Hwu (UIUC)
 - From The University of Illinois ECE 498AI class
- Some slides courtesy Massimiliano Fatica (NVIDIA)
- Many slides courtesy Kayvon Fatahalian (Stanford)

Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 2

A GPU Renders 3D Scenes

- A **Graphics Processing Unit (GPU)** accelerates rendering of 3D scenes
 - Input: description of scene
 - Output: colored pixels to be displayed on a screen
- Input:
 - Geometry (triangles), colors, lights, effects, textures
- Output:
 

Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 3

GPU Scene Complexity Defined by Standard Interfaces (DirectX and OpenGL)

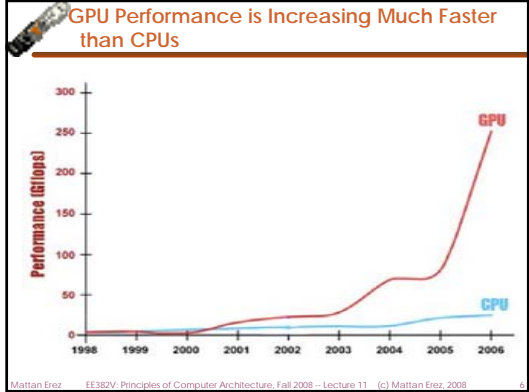
- DirectX and OpenGL define the interface between applications and the GPU
- **Geometry** describes the objects and layout
 - Triangles (vertices) describe all objects
 - Can have millions of triangles per scene
 - Can modify triangle surfaces
 - Bumps, ripples, ...
 - Lights are part of the scene geometry
- **Pixel Shaders** describe how to add color
 - Colors of triangle vertices
 - Textures (patterns)
 - How to determine color of pixels within a triangle
 - ...

Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 4

Complexity and Quality are Orders of Magnitude Better



Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 5

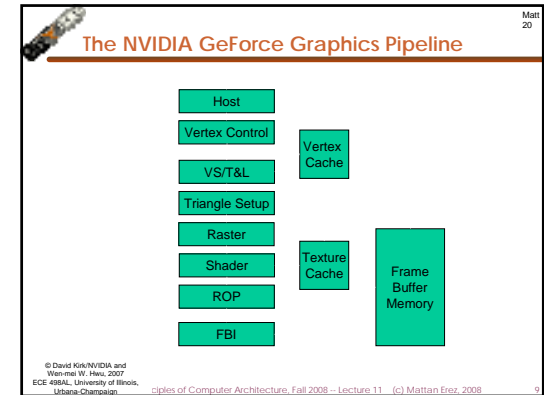


GPU and CPU Architectures are Starting to Converge

	CPUs	GPUs
1997	no explicit parallelism	not programmable
2000	explicit short vectors	emerging programmability (2001 – 2002), "infinite" DP
2003	explicit short vectors explicit threading (~2)	fully programmable explicit "infinite" DP no scatter
2006	explicit short vectors explicit threading (~4)	explicit vectors explicit threading (~16)
2009?	explicit vectors explicit threading (>16)	explicit vectors explicit threading (>16)

© David Kirk/NVIDIA and Wen-Mei Hwu, 2007
ECE 498AL, University of Illinois, Urbana-Champaign

- ### Outline
- What is a GPU?
 - Why should we care about GPUs?
 - 3D graphics pipeline
 - Programmable GPUs
-
- Many slides courtesy David Kirk (NVIDIA) and Wen-Mei Hwu (UIUC)
 - From The University of Illinois ECE 498AI class
 - Other slides courtesy Massimiliano Fatics (NVIDIA)
- © David Kirk/NVIDIA and Wen-Mei Hwu, 2007
ECE 498AL, University of Illinois, Urbana-Champaign



- ### Filtering techniques
- Point sampling:
 - pixel values are calculated by choosing one texture pixel (texel) color
 - Bilinear filtering:
 - interpolating colors from 4 neighboring texels. This gives a smoothing (if somewhat blurry) effect and makes the scene look more natural and prevents abrupt transitions between neighboring texels.
 - Trilinear filtering:
 - interpolating bilinearly filtered samples from two mip-maps. Trilinear mip-mapping prevents moving objects from displaying a distracting "sparkle" caused by abrupt transitions between mipmaps.
 - Anisotropic filtering:
 - interpolating and filtering multiple samples from one or more mip-maps to better approximate very distorted textures. Gives a sharper effect when severe perspective correction is used. Trilinear mipmapping blurs textures more.
- © David Kirk/NVIDIA and Wen-Mei Hwu, 2007
ECE 498AL, University of Illinois, Urbana-Champaign

- ### Texture Cache
-
- Stores temporally local texel values to reduce bandwidth requirements
 - Due to nature of texture filtering high degrees of efficiency are possible
 - Efficient texture caches can achieve 75% or better hit rates
 - Reduces texture (memory) bandwidth by a factor of four for bilinear filtering
- © David Kirk/NVIDIA and Wen-Mei Hwu, 2007
ECE 498AL, University of Illinois, Urbana-Champaign

- ### Pixel Shading
-
- 1999 (DirectX 7)
 - Application could select from a few simple combinations of texture and interpolated color
 - Add
 - Decal
 - Modulate
 - Next (DirectX 9)
 - Write a general program that executes for every pixel with a nearly unlimited number of interpolated inputs, texture lookups and math operations
 - Can afford to perform sophisticated lighting calculations at every pixel
- ```

pixel shader
void main() {
 tex = tex2D(Sampler0, texCoord);
 color = tex * interpColor;
}

pixel shader
void main() {
 tex = tex2D(Sampler0, texCoord);
 color = tex + interpColor;
}

pixel shader
void main() {
 tex = tex2D(Sampler0, texCoord);
 color = tex * interpColor;
}

pixel shader
void main() {
 tex = tex2D(Sampler0, texCoord);
 color = tex + interpColor;
}

```
- © David Kirk/NVIDIA and Wen-Mei Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign




## Z Buffer

- A Z buffer is a 2-D array of Z values with the same (x,y) dimensions as the color framebuffer
- Every candidate pixel from the shader has a calculated Z value along with its R,G,B,A color.
- Before writing the color, perform the Z-buffer test:
  - Read the Z value from memory
  - Compare the candidate Z to the Z from memory; if the candidate Z is NOT in front of the previous Z, discard the pixel
  - Otherwise, write the new Z value to the Z buffer and write (or blend) the new color to the color framebuffer

© David Kirk/NVIDIA and Wen-Mei Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign  
Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 19

## Summary, so far...

- Introduction to several key 3D graphics concepts:
  - Framebuffers
  - Object Representation
  - Vertex Processing
  - Lighting
  - Rasterization
  - Gouraud Interpolation
  - Texture Mapping
  - Pixel Shading
  - Alpha Blending
  - Anti-Aliasing
  - Z-Buffering



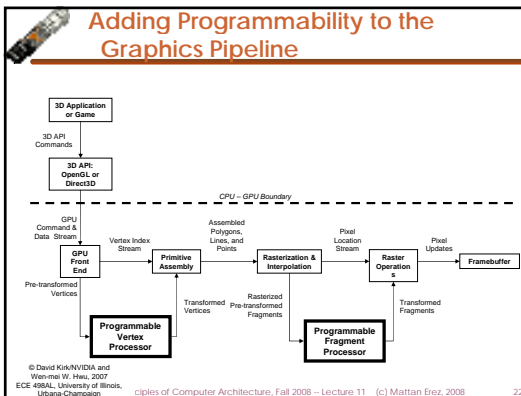
© David Kirk/NVIDIA and Wen-Mei Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign  
Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 20

## Outline

- What is a GPU?
- Why should we care about GPUs?
- 3D graphics pipeline
- Programmable GPUs

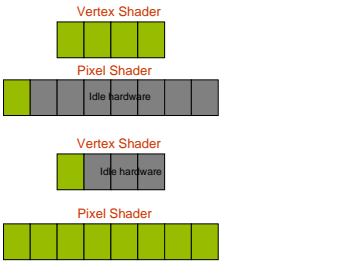
- Many slides courtesy David Kirk (NVIDIA) and Wen-Mei Hwu (UIUC)
  - From The University of Illinois ECE 498A class
- Other slides courtesy Massimiliano Faticis (NVIDIA)

Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 21



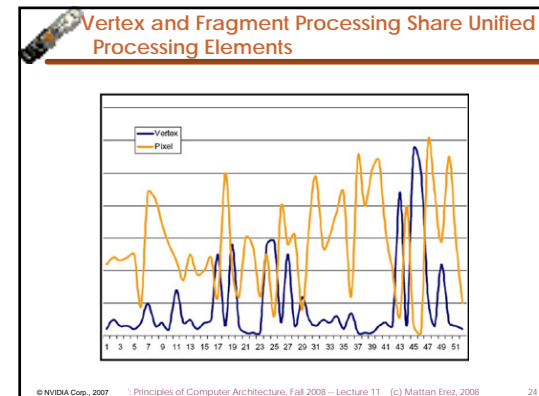
## Vertex and Fragment Processing Share Unified Processing Elements

- Load balancing HW is a problem



The diagram shows two scenarios of hardware utilization. In the 'Heavy Geometry' scenario, the workload is low (Perf = 4), meaning the vertex shader is busy while the pixel shader is idle. In the 'Heavy Pixel' scenario, the workload is high (Perf = 8), meaning the pixel shader is busy while the vertex shader is idle. This illustrates the problem of load balancing hardware resources.

© NVIDIA Corp., 2007  
Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 23



### Vertex and Fragment Processing Share Unified Processing Elements

- Load balancing SW is easier

Unified Shader

Vertex Workload

Pixel

Heavy Geometry  
Workload Perf = 11

Unified Shader

Pixel Workload

Vertex

Heavy Pixel  
Workload Perf = 11

© NVIDIA Corp., 2007 Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 24

### Vertex and Fragment Processing is Dynamically Load Balanced

Less Geometry

More Geometry

Unified Shader Usage

High pixel shader use  
Low vertex shader use

Balanced use of pixel shader and vertex shader

© NVIDIA Corp., 2007 Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 25

### Make the Compute Core The Focus of the Architecture

- The future of GPUs is programmable processing
- Processors execute competing threads
- So build the architecture around the processor
- Alternative operating mode specifically for computing

ECE 498AL, University of Illinois, Urbana-Champaign Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 27

### The NVIDIA GeForce Graphics Pipeline

© David Kirk/NVIDIA and Willem W. Huis, 2007 ECE 498AL, University of Illinois, Urbana-Champaign Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 28

### Another View of the 3D Graphics Pipeline

Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 29

### The NVIDIA GeForce Graphics Pipeline

Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 30

### Stream Execution Model

- Data parallel *streams* of data
- Processing *kernels*
  - Unit of Execution is processing of one stream element in one kernel – defined as a *thread*

Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 31

### Stream Execution Model

- Can partition the streams into chunks
  - Streams are very long and elements are independent
  - Chunks are called *strips* or *blocks*

- Unit of Execution is processing one block of data by one kernel – defined as a *thread block*

Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 32

## From Shader Code to a Teraflop: How Shader Cores Work

Kayvon Fatahalian  
Stanford University

Kayvon Fatahalian@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Mattan Erez, 2008 33

### This talk

- Three key concepts behind how modern architectures run “shader” code
- Knowing these concepts will help you:
  - Understand space of GPU shader core (and throughput CPU processing core) designs
  - Optimize shaders/compute kernels
  - Establish intuition: what workloads might benefit from the design of these architectures?

Kayvon Fatahalian@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Kayvon Fatahalian, 2008 34

### What’s in a GPU?

Heterogeneous chip multi-processor (highly tuned for graphics)

Kayvon Fatahalian@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Kayvon Fatahalian, 2008 35

### A diffuse reflectance shader

```

sampler mySamp;
Texture2D<float3>myTex;
float3 lightDir;

float4 diffuseShader(float3 norm, float2 uv)
{
 float3 kd;
 kd = myTex.Sample(mySamp, uv);
 kd *= clamp(dot(lightDir, norm), 0.0, 1.0);
 return float4(kd, 1.0);
}

```

Independent – data parallelism

Kayvon Fatahalian@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (c) Kayvon Fatahalian, 2008 36

### Compile shader

1 unshaded fragment input record

```

sampler mySamp;
Texture2D(float3)myTex;
float3 lightDir;

float4 diffuseShader(float3 norm, float2 uv)
{
 float3 kd;
 kd = myTex.Sample(mySamp, uv);
 kd *= clamp (dot(lightDir, norm), 0.0, 1.0);
 return float4(kd, 1.0);
}

```

```

<diffuseShader>;
sample r0, v4, t0, s0
mul r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul o0, r0, r3
mul o1, r1, r3
mul o2, r2, r3
mov o3, 1(1.0)

```

1 shaded fragment output record

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 37

### Execute shader

```

<diffuseShader>;
sample r0, v4, t0, s0
mul r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul o0, r0, r3
mul o1, r1, r3
mul o2, r2, r3
mov o3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 38

### Execute shader

```

<diffuseShader>;
sample r0, v4, t0, s0
mul r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul o0, r0, r3
mul o1, r1, r3
mul o2, r2, r3
mov o3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 39

### Execute shader

```

<diffuseShader>;
sample r0, v4, t0, s0
mul r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul o0, r0, r3
mul o1, r1, r3
mul o2, r2, r3
mov o3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 40

### Execute shader

```

<diffuseShader>;
sample r0, v4, t0, s0
mul r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul o0, r0, r3
mul o1, r1, r3
mul o2, r2, r3
mov o3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 41

### Execute shader

```

<diffuseShader>;
sample r0, v4, t0, s0
mul r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul o0, r0, r3
mul o1, r1, r3
mul o2, r2, r3
mov o3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 42

### Execute shader

```

diffuseShader:
sample r3, v4, tex, s8
mul r3, v0, c0[0]
madd r3, v1, c0[1], r3
madd r3, v2, c0[2], r3
clamp r3, r3, [0.0], [1.0]
mul s0, r0, r3
mul s1, r1, r3
mul s2, r2, r3
mov s3, [1.0]

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalian, 2008 43

### CPU-"style" cores

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalian, 2008 44

### Slimming down

Idea #1:  
Remove components that help a single instruction stream run fast

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalian, 2008 45

### Two cores (two fragments in parallel)

fragment 1

```

fragment 1:
sample r3, v4, tex, s8
mul r3, v0, c0[0]
madd r3, v1, c0[1], r3
madd r3, v2, c0[2], r3
clamp r3, r3, [0.0], [1.0]
mul s0, r0, r3
mul s1, r1, r3
mul s2, r2, r3
mov s3, [1.0]

```

fragment 2

```

fragment 2:
sample r3, v4, tex, s8
mul r3, v0, c0[0]
madd r3, v1, c0[1], r3
madd r3, v2, c0[2], r3
clamp r3, r3, [0.0], [1.0]
mul s0, r0, r3
mul s1, r1, r3
mul s2, r2, r3
mov s3, [1.0]

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalian, 2008 46

### Four cores (four fragments in parallel)

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalian, 2008 47

### Sixteen cores (sixteen fragments in parallel)

16 cores = 16 simultaneous instruction streams

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalian, 2008 48



### Instruction stream coherence

But... many fragments should be able to share an instruction stream!

```

<diffuseShader>:
sample r0, v4, t0, s0
mul r3, v0, c0[0]
madd r3, v1, c0[1], r3
madd r3, v2, c0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul r0, r0, r3
mul r1, r1, r3
mul r2, r2, r3
mov r3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (©Kayvon Fatahalla, 2008) 49

### Recall: simple processing core

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (©Kayvon Fatahalla, 2008) 50

### Add ALUs

Idea #2:  
Amortize cost/complexity of managing an instruction stream across many ALUs

SIMD processing

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (©Kayvon Fatahalla, 2008) 51

### Modifying the shader

```

<diffuseShader>:
sample r0, v4, t0, s0
mul r3, v0, c0[0]
madd r3, v1, c0[1], r3
madd r3, v2, c0[2], r3
clamp r3, r3, 1(0.0), 1(1.0)
mul r0, r0, r3
mul r1, r1, r3
mul r2, r2, r3
mov r3, 1(1.0)

```

Original compiled shader:  
Processes one fragment using scalar ops on scalar registers

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (©Kayvon Fatahalla, 2008) 52

### Modifying the shader

```

<VECR_diffuseShader>:
VECR_sample vec_r0, vec_v4, t0, vec_s0
VECR_mul vec_r3, vec_v0, c0[0]
VECR_madd vec_r3, vec_v1, c0[1], vec_r3
VECR_madd vec_r3, vec_v2, c0[2], vec_r3
VECR_clamp vec_r3, vec_r3, 1(0.0), 1(1.0)
VECR_mul vec_r0, vec_r0, vec_r3
VECR_mul vec_r1, vec_r1, vec_r3
VECR_mul vec_r2, vec_r2, vec_r3
VECR_mov vec_r3, 1(1.0)

```

New compiled shader:  
Processes 8 fragments using vector ops on vector registers

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (©Kayvon Fatahalla, 2008) 53

### Modifying the shader

```

<VECR_diffuseShader>:
VECR_sample vec_r0, vec_v4, t0, vec_s0
VECR_mul vec_r3, vec_v0, c0[0]
VECR_madd vec_r3, vec_v1, c0[1], vec_r3
VECR_madd vec_r3, vec_v2, c0[2], vec_r3
VECR_clamp vec_r3, vec_r3, 1(0.0), 1(1.0)
VECR_mul vec_r0, vec_r0, vec_r3
VECR_mul vec_r1, vec_r1, vec_r3
VECR_mul vec_r2, vec_r2, vec_r3
VECR_mov vec_r3, 1(1.0)

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11 (©Kayvon Fatahalla, 2008) 54

### 128 fragments in parallel

16 cores = 128 ALUs  
= 16 simultaneous instruction streams

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 11 (c)Kayvon Fatahalla, 2008 56

### 128 [ primitives, vertices, fragments ] in parallel

Vertices/Fragments  
primitives  
CUDA threads  
OpenGL work items  
compute shader threads

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 11 (c)Kayvon Fatahalla, 2008 56

### But what about branches?

```

<unconditional
shader code>
if (x > 0) {
 y = pow(x, exp);
 y *= Ks;
 refl = y + Ka;
} else {
 x = 0;
 refl = Ka;
}
<resume unconditional
shader code>

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 11 (c)Kayvon Fatahalla, 2008 57

### But what about branches?

```

<unconditional
shader code>
if (x > 0) {
 y = pow(x, exp);
 y *= Ks;
 refl = y + Ka;
} else {
 x = 0;
 refl = Ka;
}
<resume unconditional
shader code>

```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 11 (c)Kayvon Fatahalla, 2008 58

### But what about branches?

```

<unconditional
shader code>
if (x > 0) {
 y = pow(x, exp);
 y *= Ks;
 refl = y + Ka;
} else {
 x = 0;
 refl = Ka;
}
<resume unconditional
shader code>

```

Not all ALUs do useful work!  
Worst case: 1/8 performance

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 11 (c)Kayvon Fatahalla, 2008 59

### But what about branches?

```

<unconditional
shader code>
if (x > 0) {
 y = pow(x, exp);
 y *= Ks;
 refl = y + Ka;
} else {
 x = 0;
 refl = Ka;
}
<resume unconditional
shader code>


```

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 11 (c)Kayvon Fatahalla, 2008 60

### Clarification

SIMD processing does not imply SIMD instructions

- Option 1: Explicit vector instructions
  - Intel/AMD x86 SSE, Intel Larrabee
- Option 2: Scalar instructions, implicit HW vectorization
  - HW determines instruction stream sharing across ALUs (amount of sharing hidden from software)
  - NVIDIA GeForce ("SIMT" warps), ATIRadeon architectures



In practice: 16 to 64 fragments share an instruction stream

Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 61, 61

## Stalls!

Stalls occur when a core cannot run the next instruction because of a dependency on a previous operation.

Texture access latency = 100's to 1000's of cycles

We've removed the fancy caches and logic that helps avoid stalls.

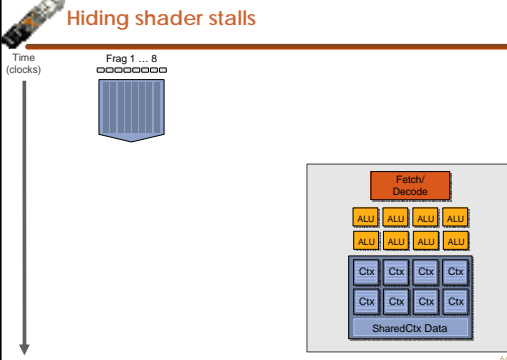
Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 62, 62

But we have LOTS of independent fragments.

**Idea #3:**  
Interleave processing of many fragments on a single core to avoid stalls caused by high latency operations.

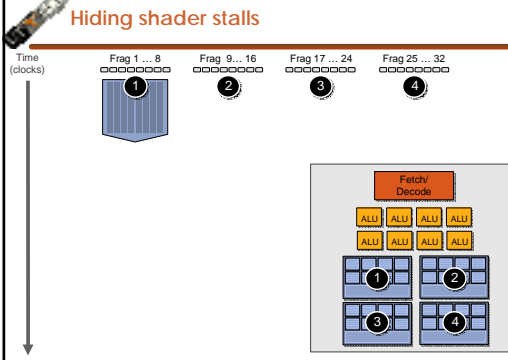
Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 63, 63

### Hiding shader stalls



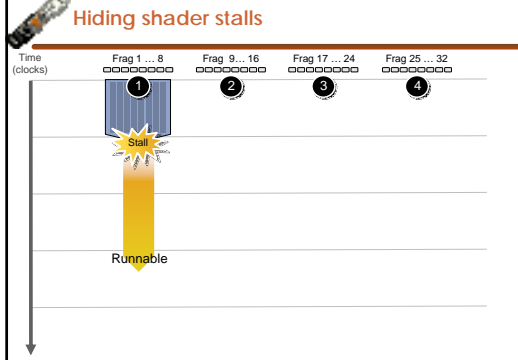
Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 64, 64

### Hiding shader stalls

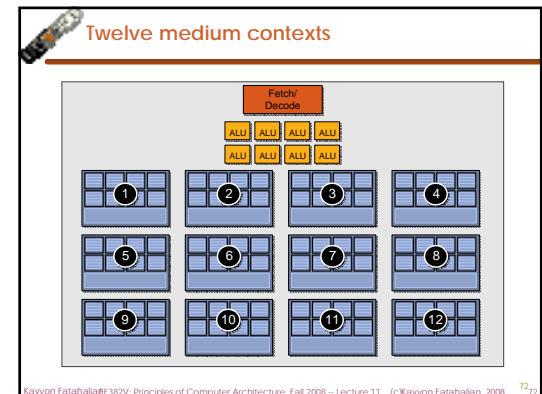
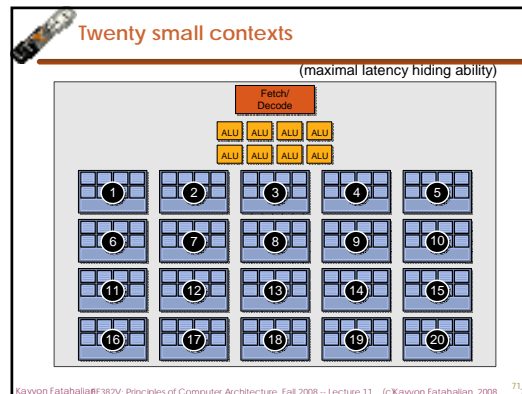
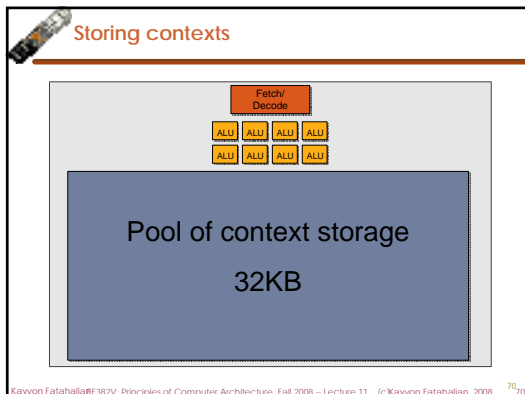
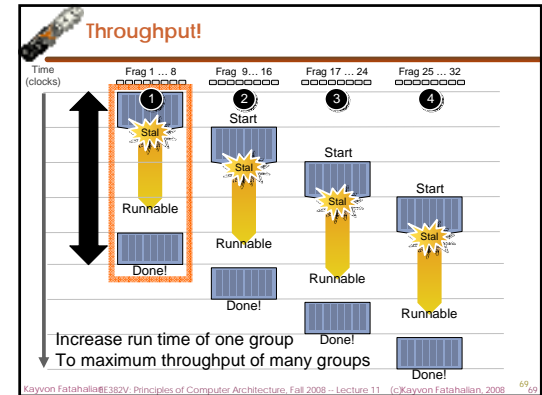
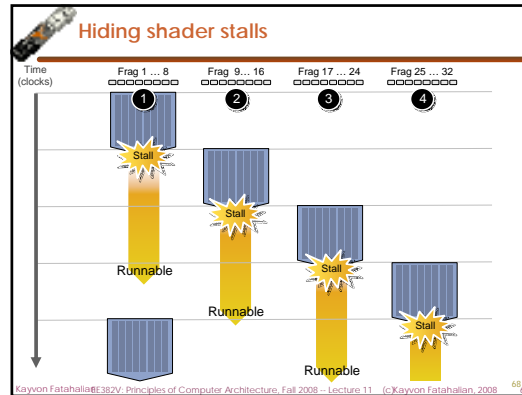
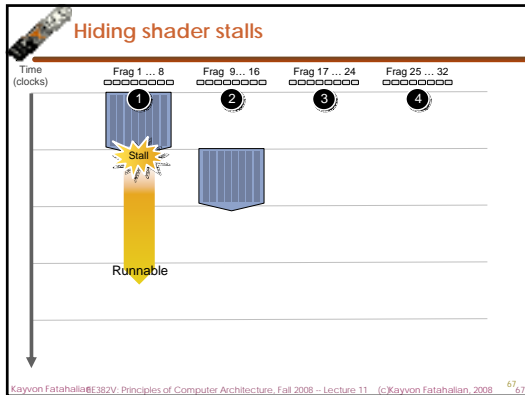


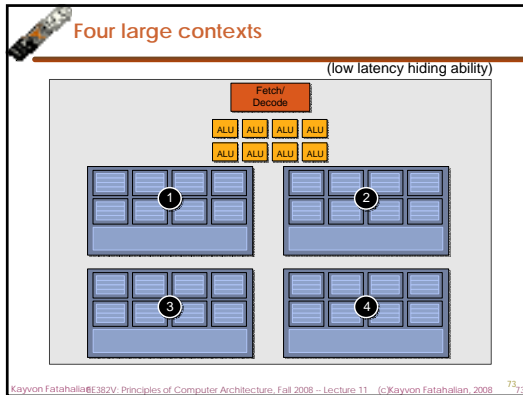
Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 65, 65

### Hiding shader stalls



Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 -- Lecture 11 (c)Kayvon Fatahalla, 2008 66, 66





- ### Summary: three key ideas
1. Use many “slimmed down cores” to run in parallel
  2. Pack cores full of ALUs (by sharing instruction stream across groups of fragments)
    - Option 1: Explicit SIMD vector instructions
    - Option 2: Implicit sharing managed by hardware
  3. Avoid latency stalls by interleaving execution of many groups of fragments
    - When one group stalls, work on another group
- Kayvon Fatahalla@EE382V, Principles of Computer Architecture, Fall 2008 – Lecture 11. ©Kayvon Fatahalla, 2008. 74