

# **CS 380 - GPU and GPGPU Programming**

## **Lecture 5: GPU Architecture 3**

Markus Hadwiger, KAUST

# Reading Assignment #3 (until Sep 21)



## Read (required):

- Programming Massively Parallel Processors book, Chapter 1 (*Introduction*)
- Programming Massively Parallel Processors book (2<sup>nd</sup> edition), Appendix B (*GPU Compute Capabilities*)
- OpenGL 4 Shading Language Cookbook, Chapter 2

## Read (optional):

- OpenGL 4 Shading Language Cookbook, Chapter 1
- GLSL (orange) book, Chapter 7 (OpenGL Shading Language API)

# OpenGL Tutorial



With Peter Rautek

Tuesday, Sep 15, 13:00 – 15:00

Come with your conceptual or coding questions!

# More Motivational Examples



## Doom (2016)

[http://www.adriancourreges.com/blog/2016/09/09/  
doom-2016-graphics-study/](http://www.adriancourreges.com/blog/2016/09/09/doom-2016-graphics-study/)

## Doom Eternal

[https://simoncoenen.com/blog/programming/graphics/  
DoomEternalStudy.html](https://simoncoenen.com/blog/programming/graphics/DoomEternalStudy.html)

## Unreal Engine 5

[https://www.unrealengine.com/en-US/blog/  
a-first-look-at-unreal-engine-5](https://www.unrealengine.com/en-US/blog/a-first-look-at-unreal-engine-5)

# 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, but no explicit parallelism

# 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  
clmp  r3, r3, l(0.0), l(1.0)  
mul   o0, r0, r3  
mul   o1, r1, r3  
mul   o2, r2, r3  
mov   o3, l(1.0)
```



1 shaded fragment output record



# Per-Pixel(Fragment) Lighting

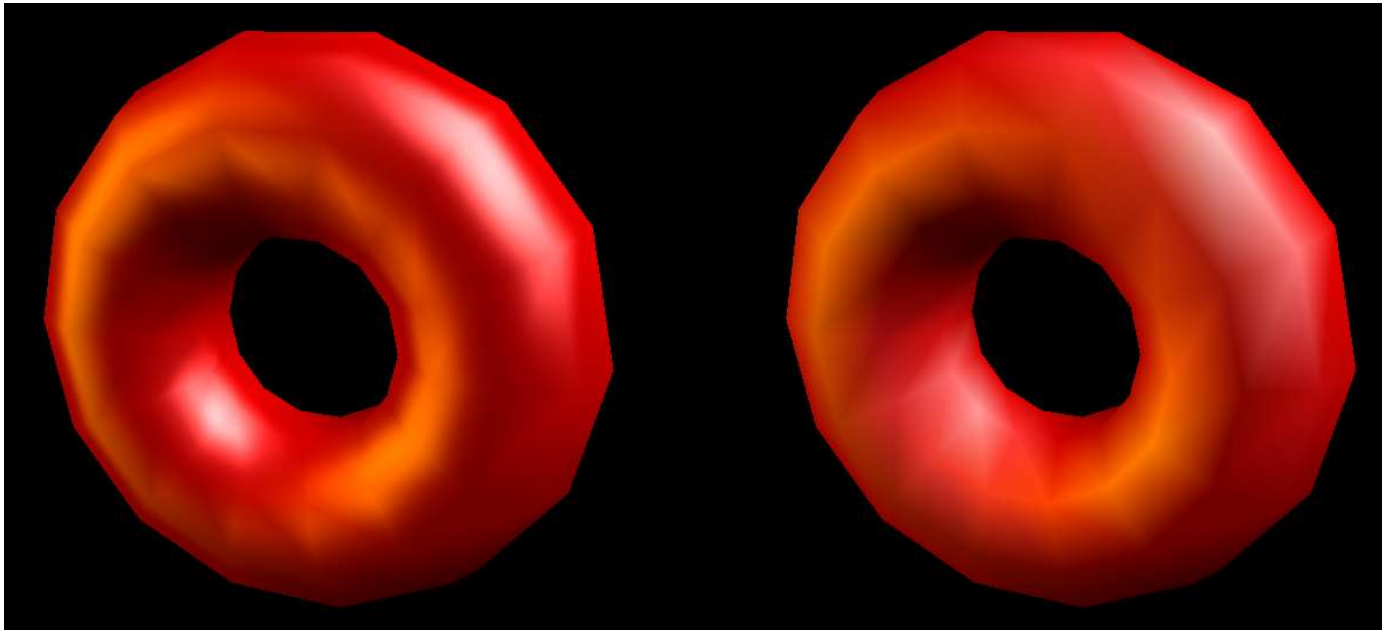


Simulating smooth surfaces by calculating illumination for each fragment

Example: specular highlights (Phong illumination/shading)

Phong shading:  
per-fragment evaluation

Gouraud shading:  
linear interpolation from vertices





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

Kayvon Fatahalian  
Stanford University



# Part 1: throughput processing

---

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

# Where this is going...

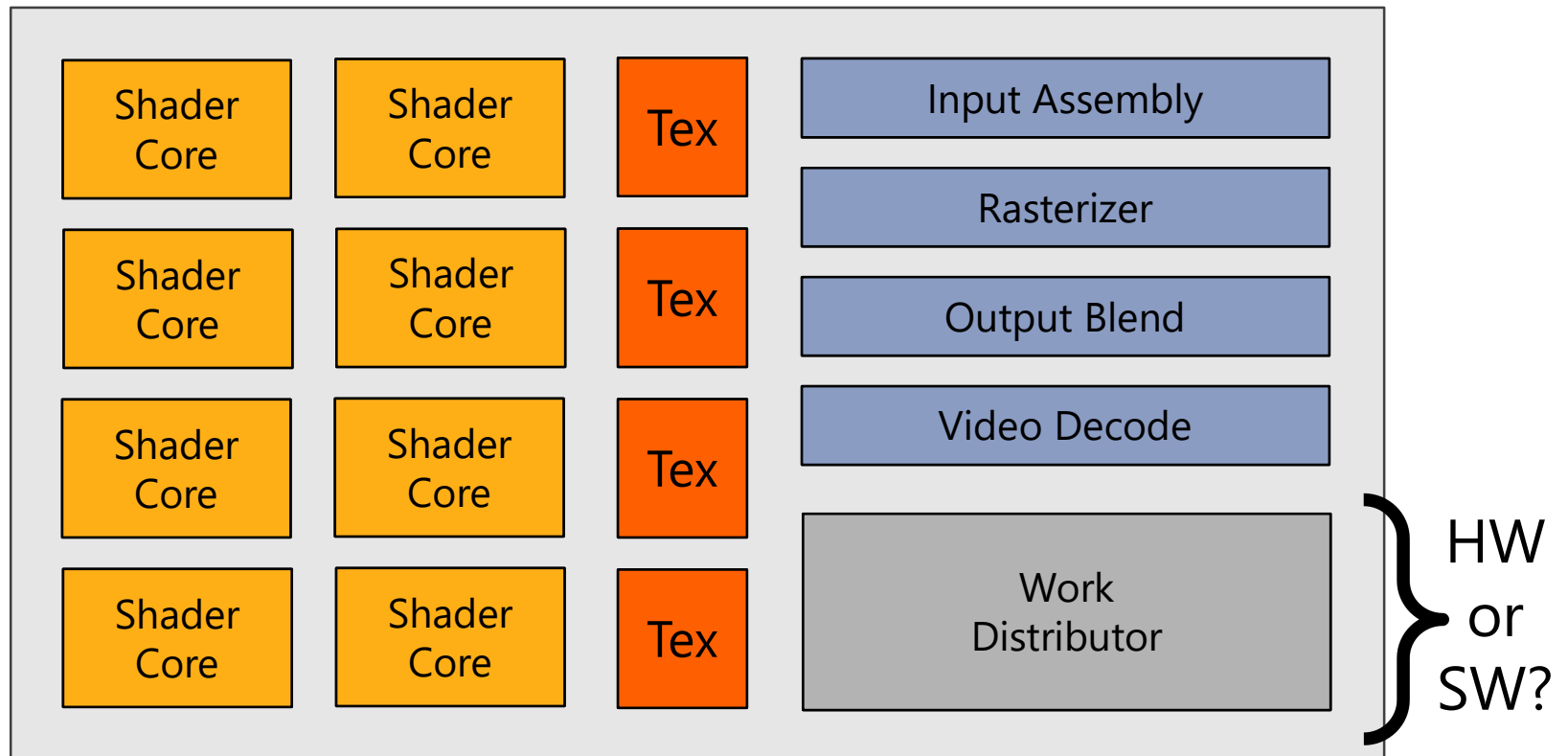


## **Summary: three key ideas for high-throughput execution**

- 1. Use many “slimmed down cores,” run them in parallel**
- 2. Pack cores full of ALUs (by sharing instruction stream overhead 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**

# What's in a GPU?

---



Heterogeneous chip multi-processor (highly tuned for graphics)

# 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, but no explicit parallelism

# 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);  
}
```



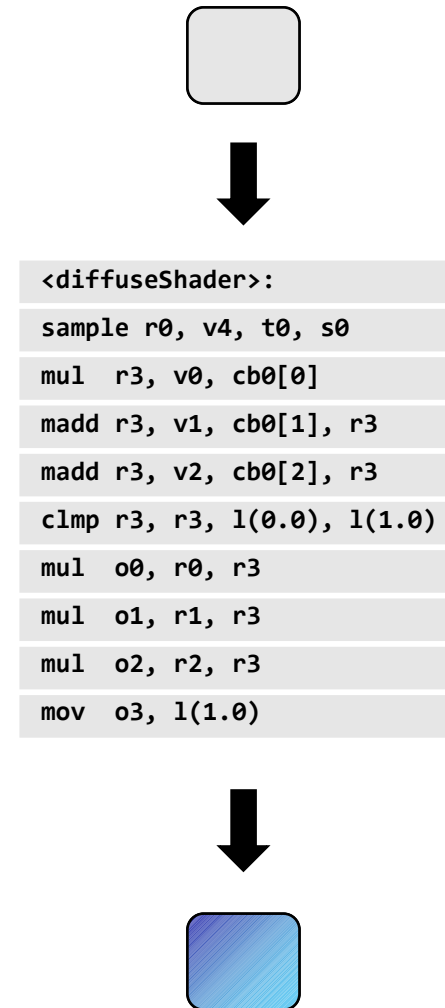
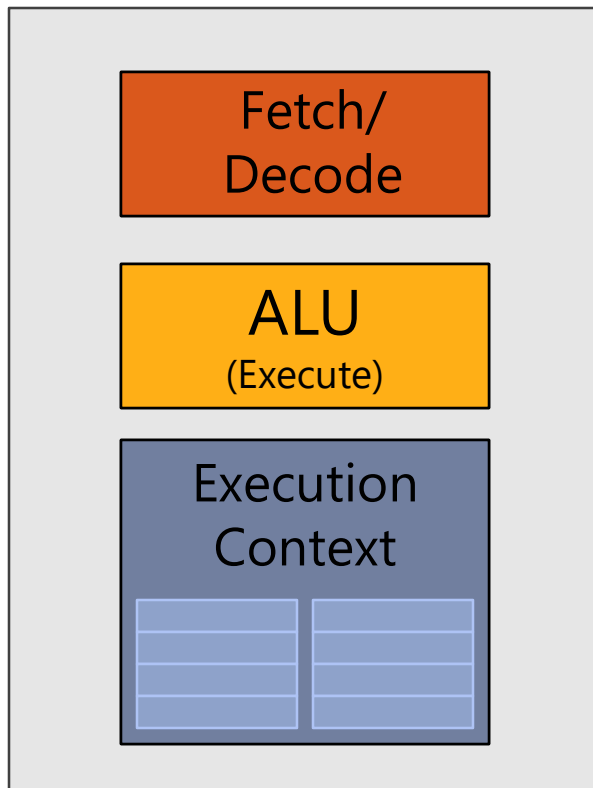
```
<diffuseShader>:  
sample r0, v4, t0, s0  
mul   r3, v0, cb0[0]  
madd  r3, v1, cb0[1], r3  
madd  r3, v2, cb0[2], r3  
clmp  r3, r3, l(0.0), l(1.0)  
mul   o0, r0, r3  
mul   o1, r1, r3  
mul   o2, r2, r3  
mov   o3, l(1.0)
```



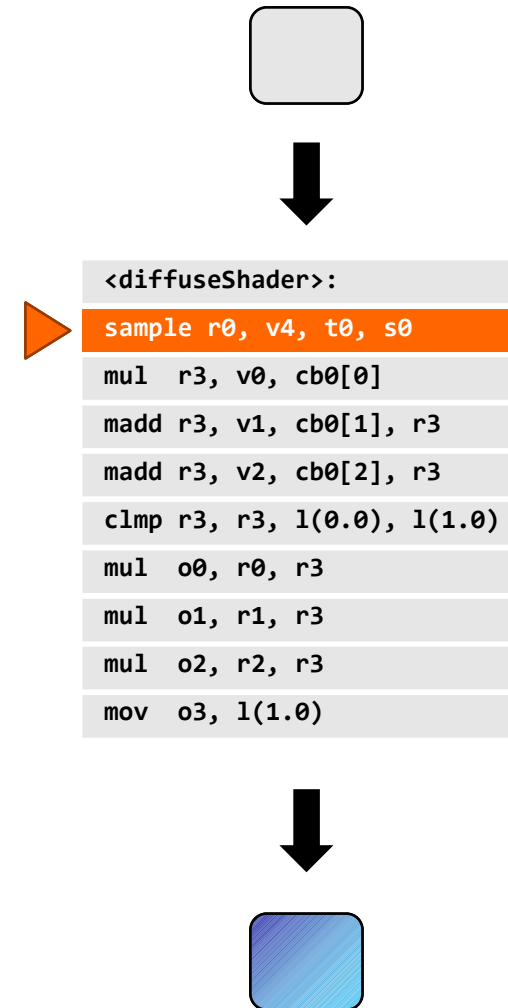
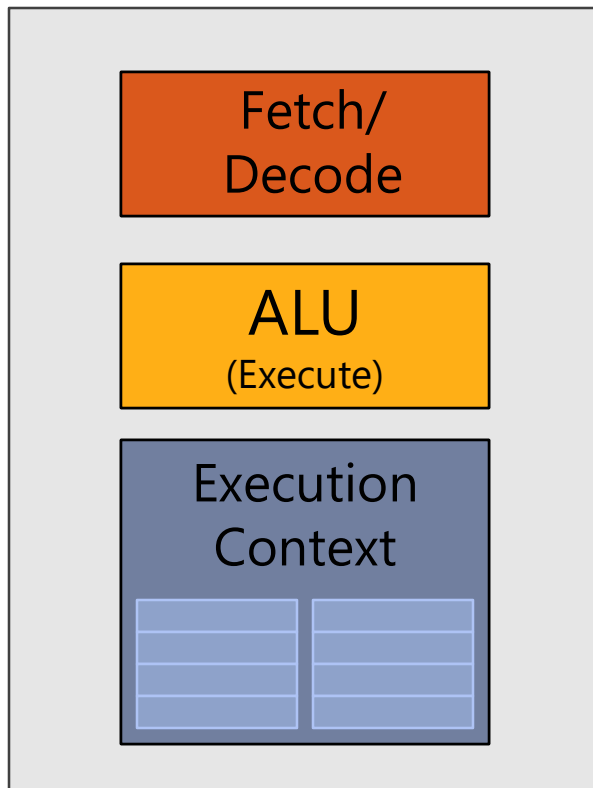
1 shaded fragment output record



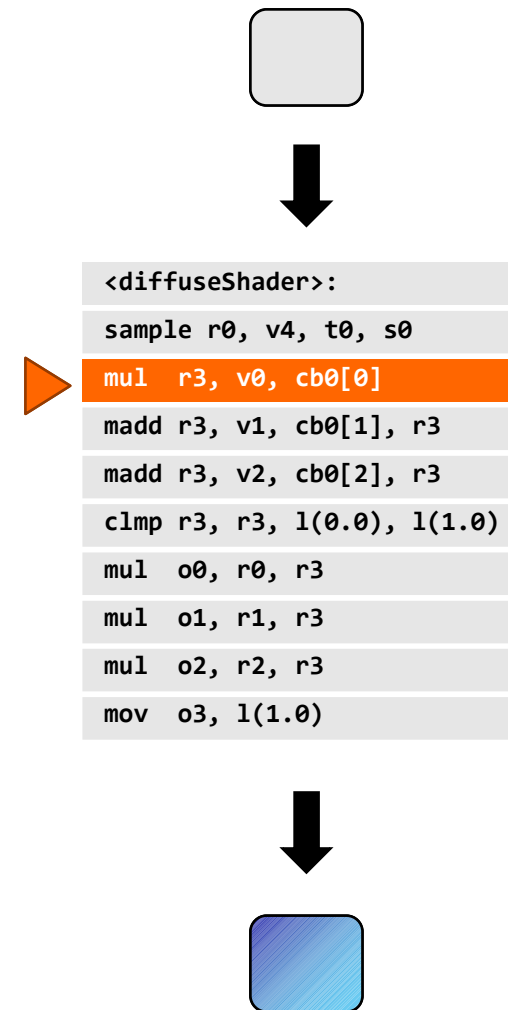
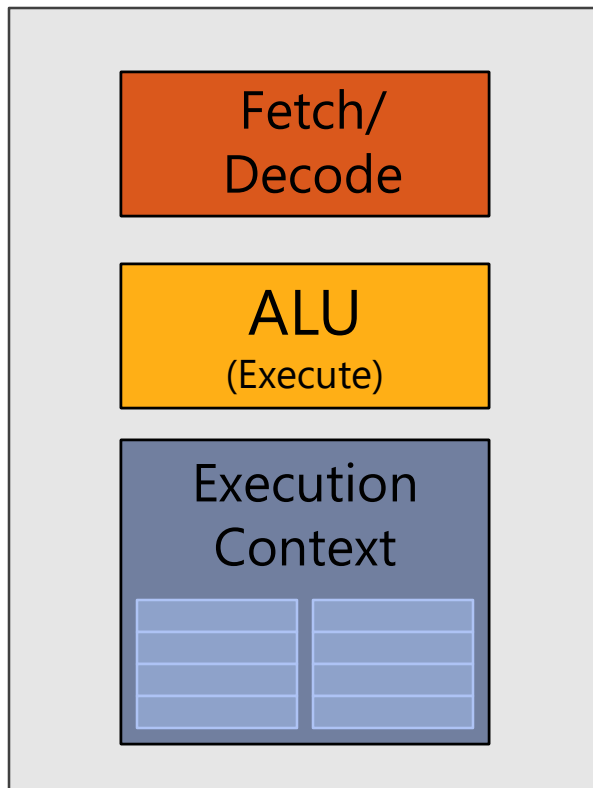
# Execute shader



# Execute shader

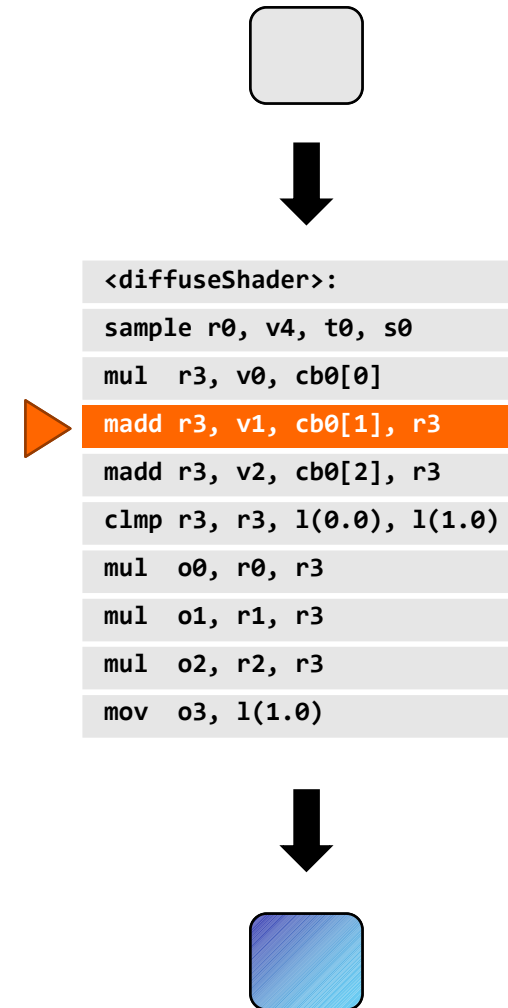
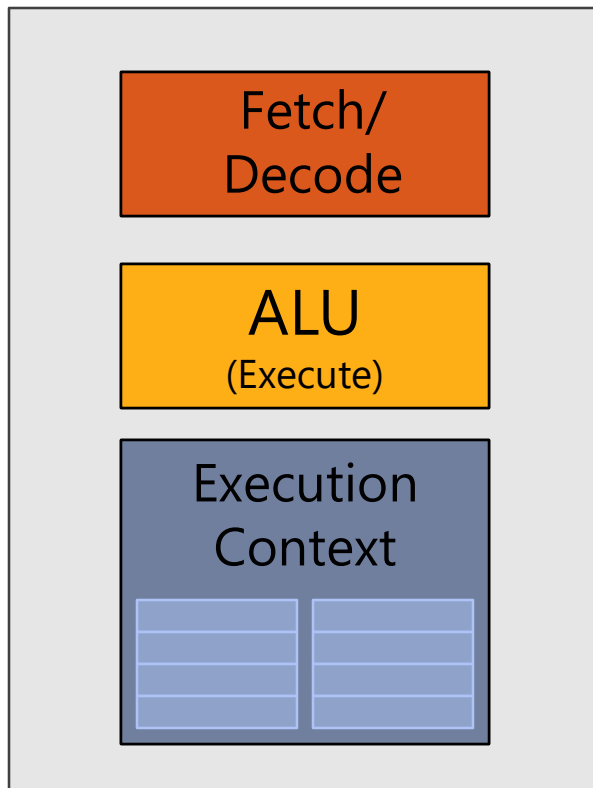


# Execute shader

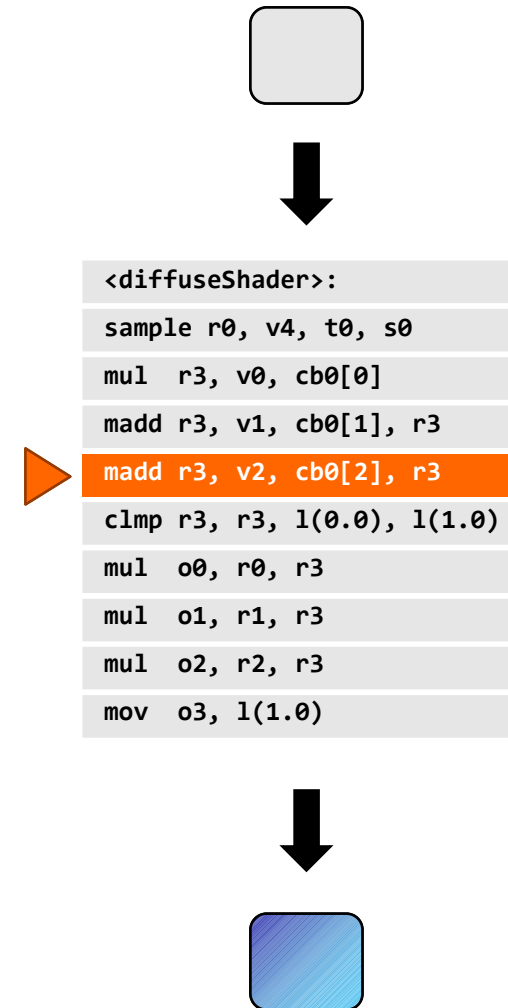
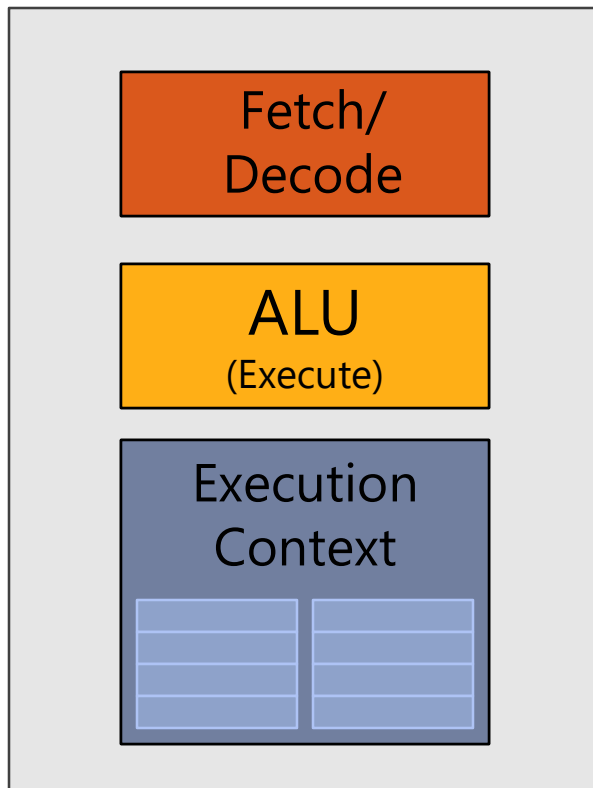




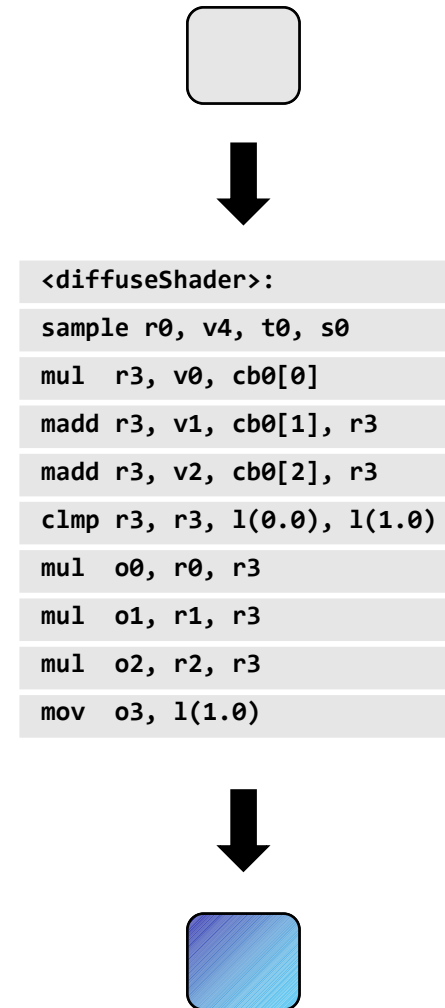
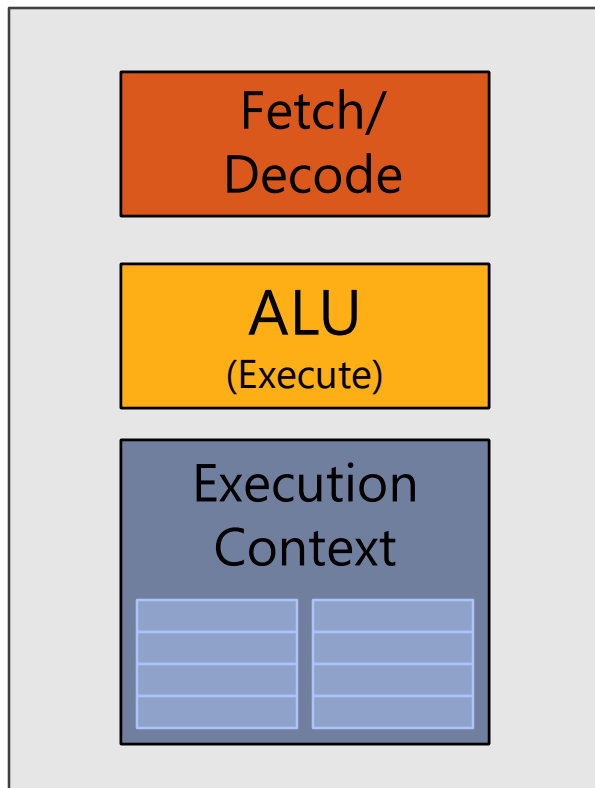
# Execute shader



# Execute shader

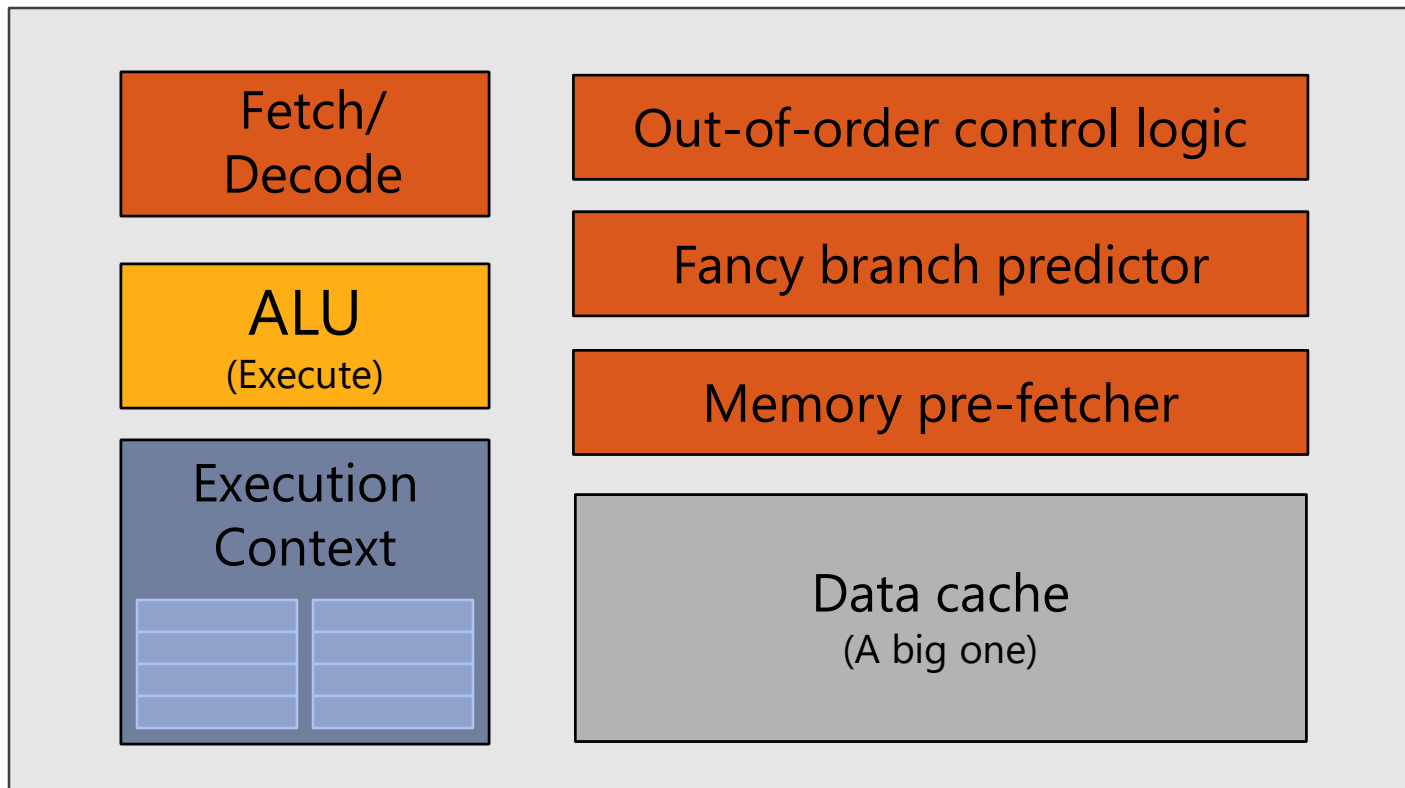


# Execute shader



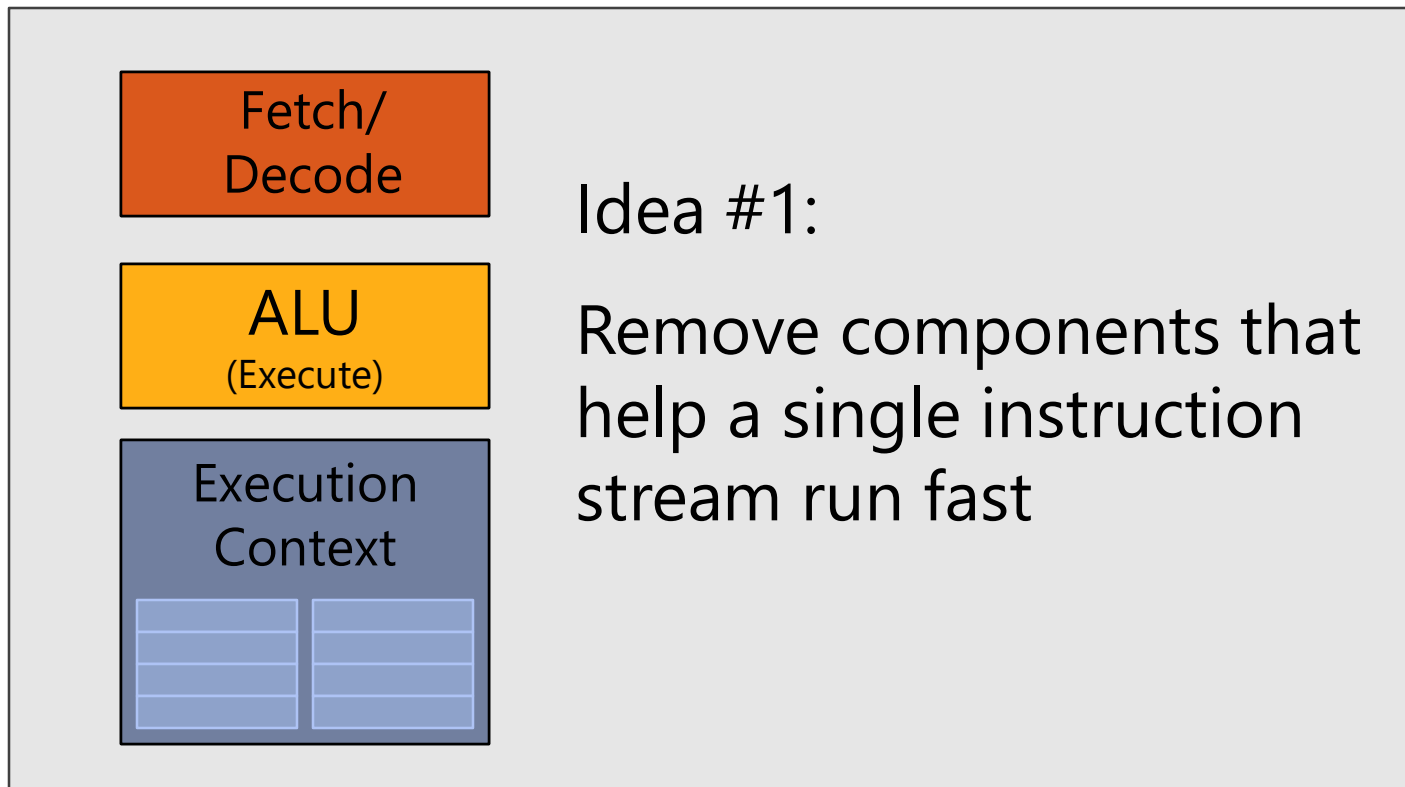
# CPU-"style" cores

---



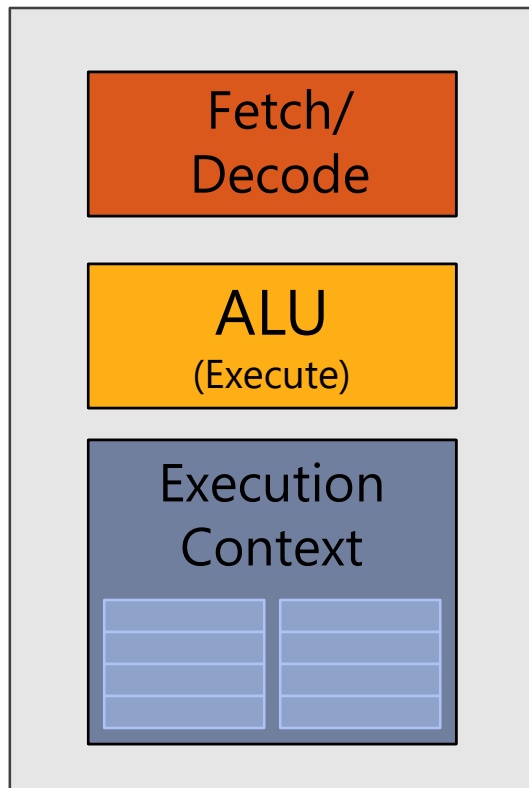
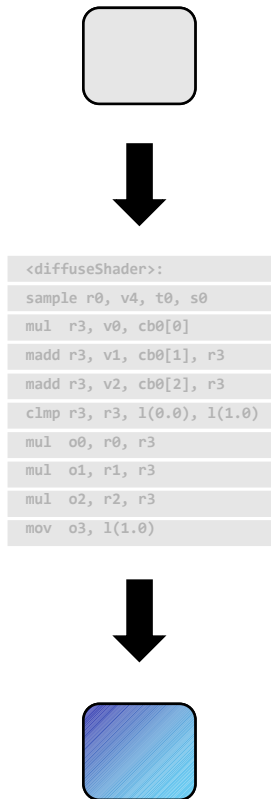
# Slimming down

---

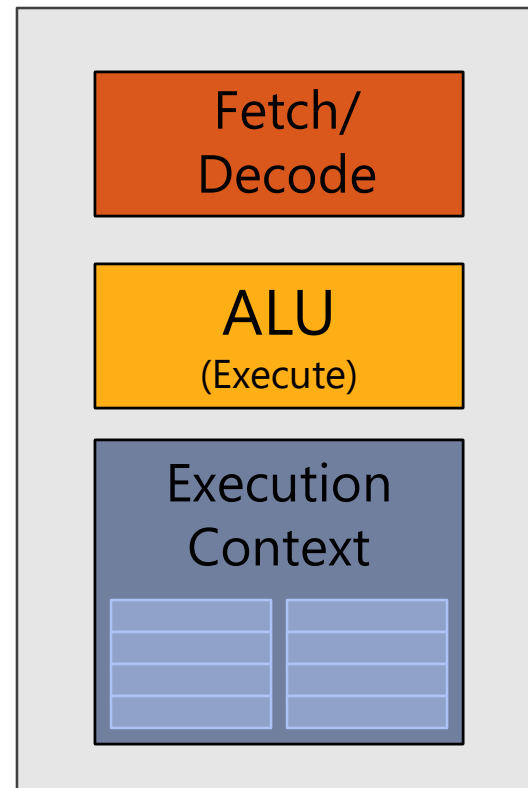
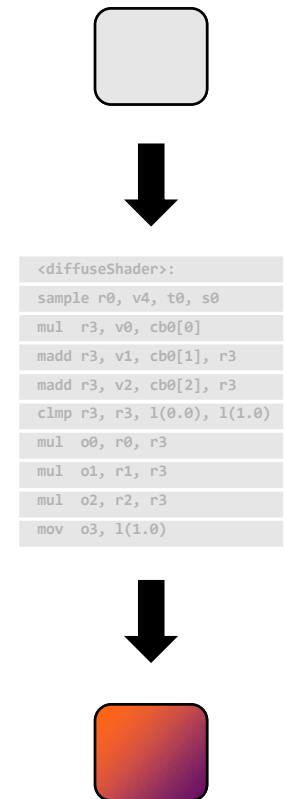


# Two cores (two fragments in parallel)

fragment 1

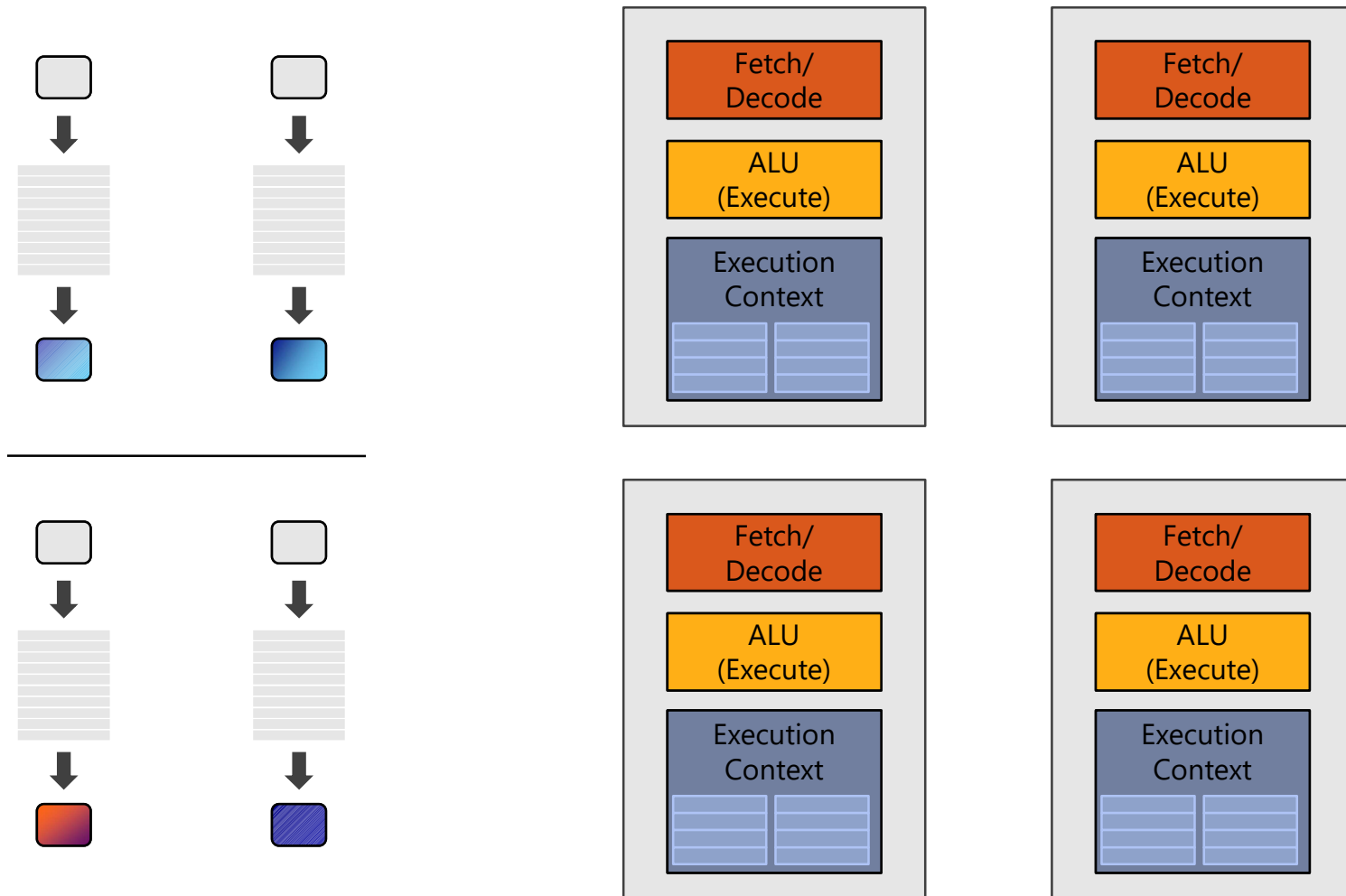


fragment 2



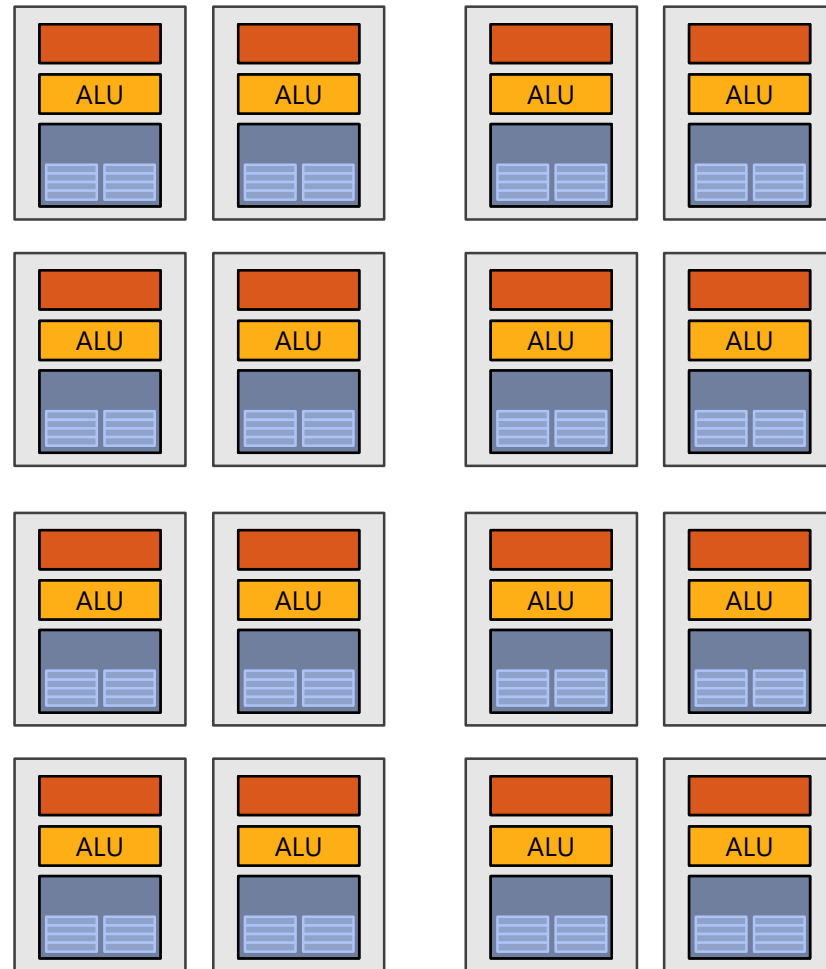
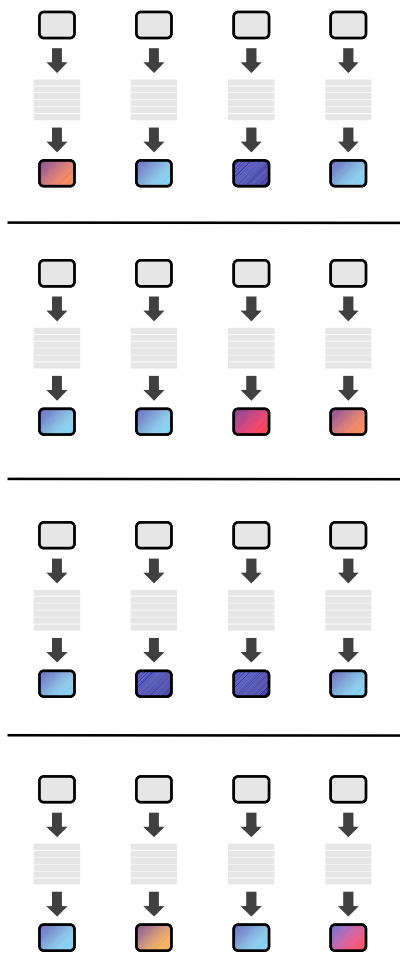
# Four cores (four fragments in parallel)

---



# Sixteen cores (sixteen fragments in parallel)

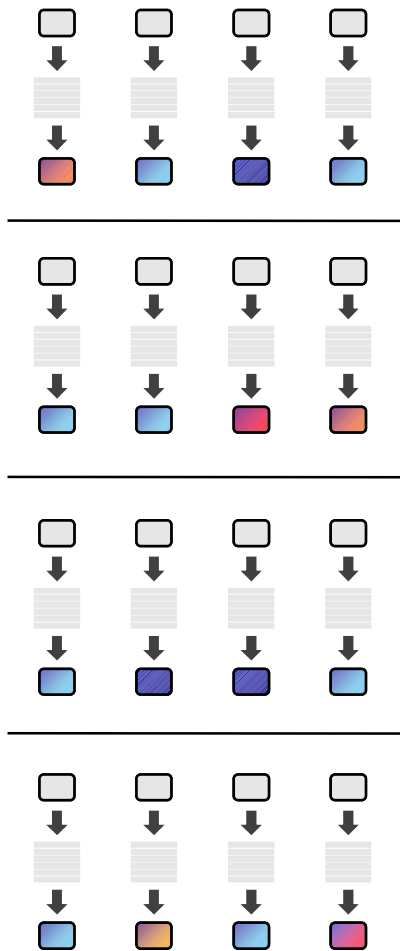
---



16 cores = 16 simultaneous instruction streams



# Instruction stream sharing

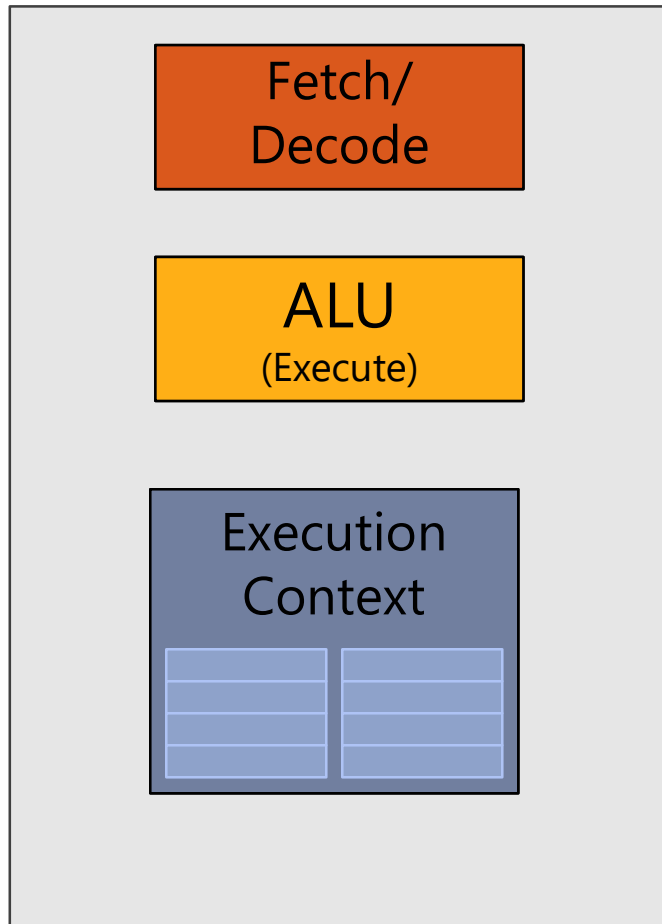


But... many fragments should be able to share an instruction stream!

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

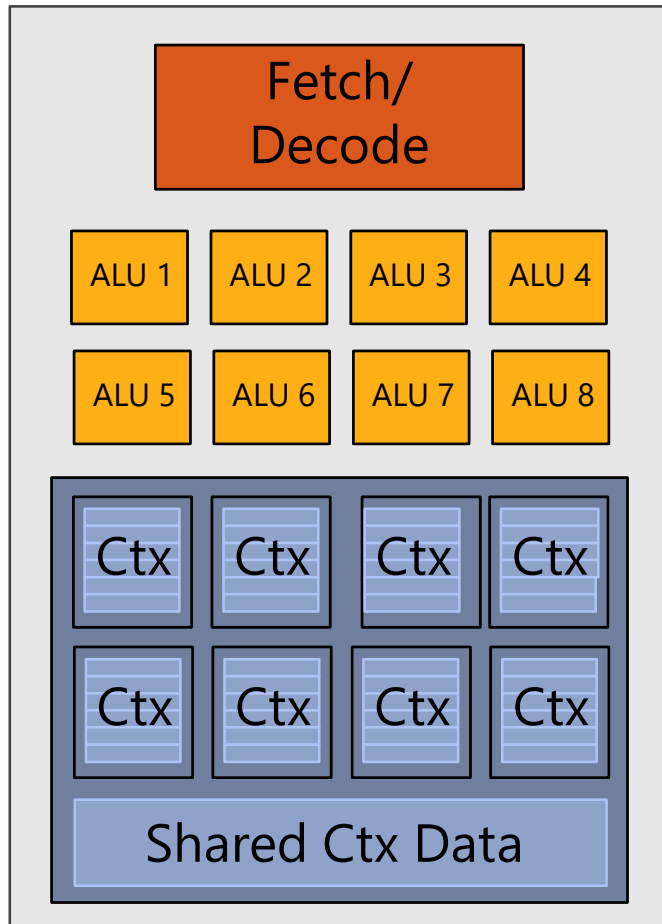
# Recall: simple processing core

---



# Add ALUs

---



Idea #2:

Amortize cost/complexity of managing an instruction stream across many ALUs

SIMD processing

(or SIMT, SPMD)

Thank you.