Dieter Schmalstieg

Shader Programming

# Today's Agenda

- The GPU execution model

- How to write shader programs

# Graphics Application Programmer Interface

- Hardware-vendor independent interface
  - Interface is hardware independent
  - But implementation is hardware dependent
- Defines
  - Abstract rendering device
  - Set of functions to operate the device

# API and Vendor Overview

Graphics API

- DirectX (Microsoft)
- OpenGL
- OpenGL ES
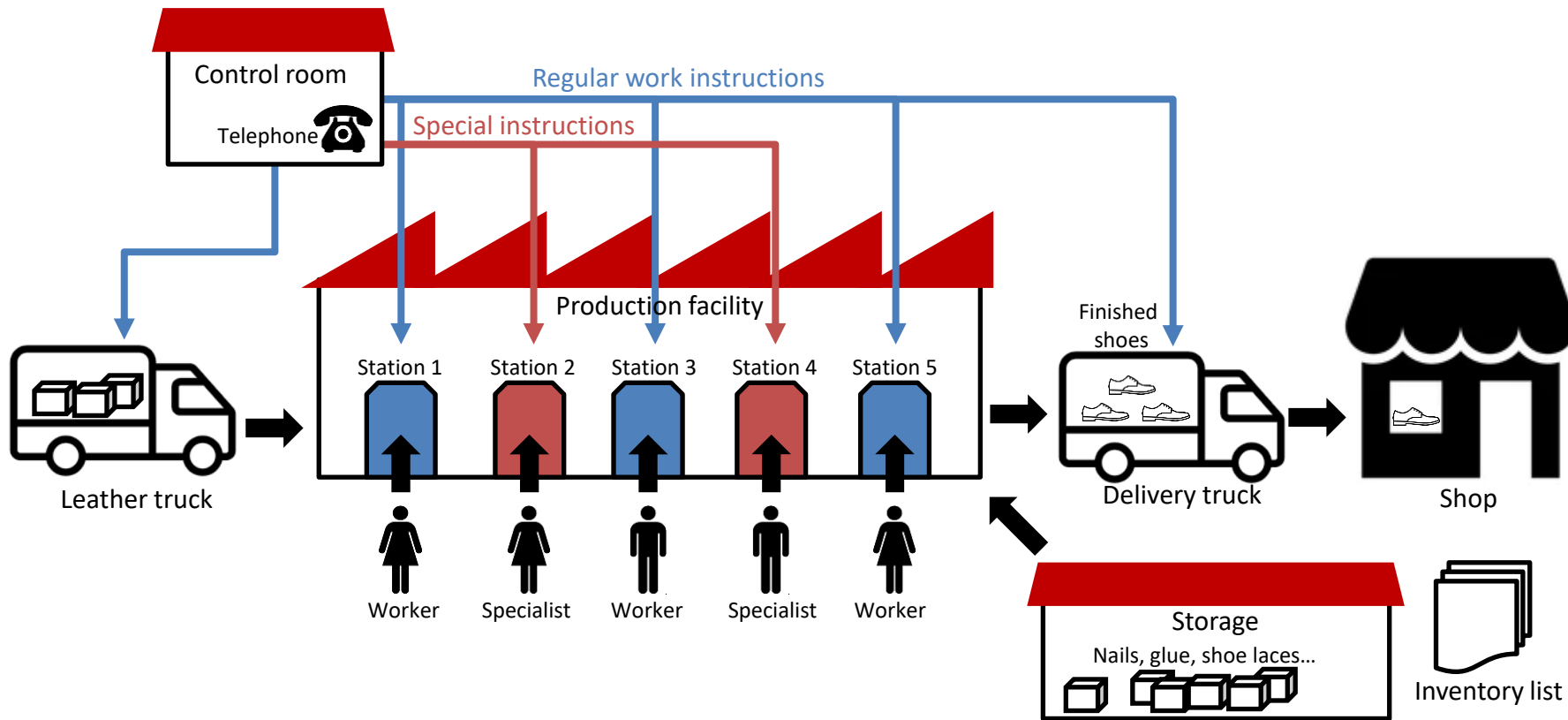- Vulkan (Khronos group)
- Metal (Apple)

GPU Hardware Vendors

- Intel (Iris, $X^e$)
- NVIDIA (GeForce)
- AMD (Radeon)
- Qualcomm (Adreno)
- Imagination (PowerVR)
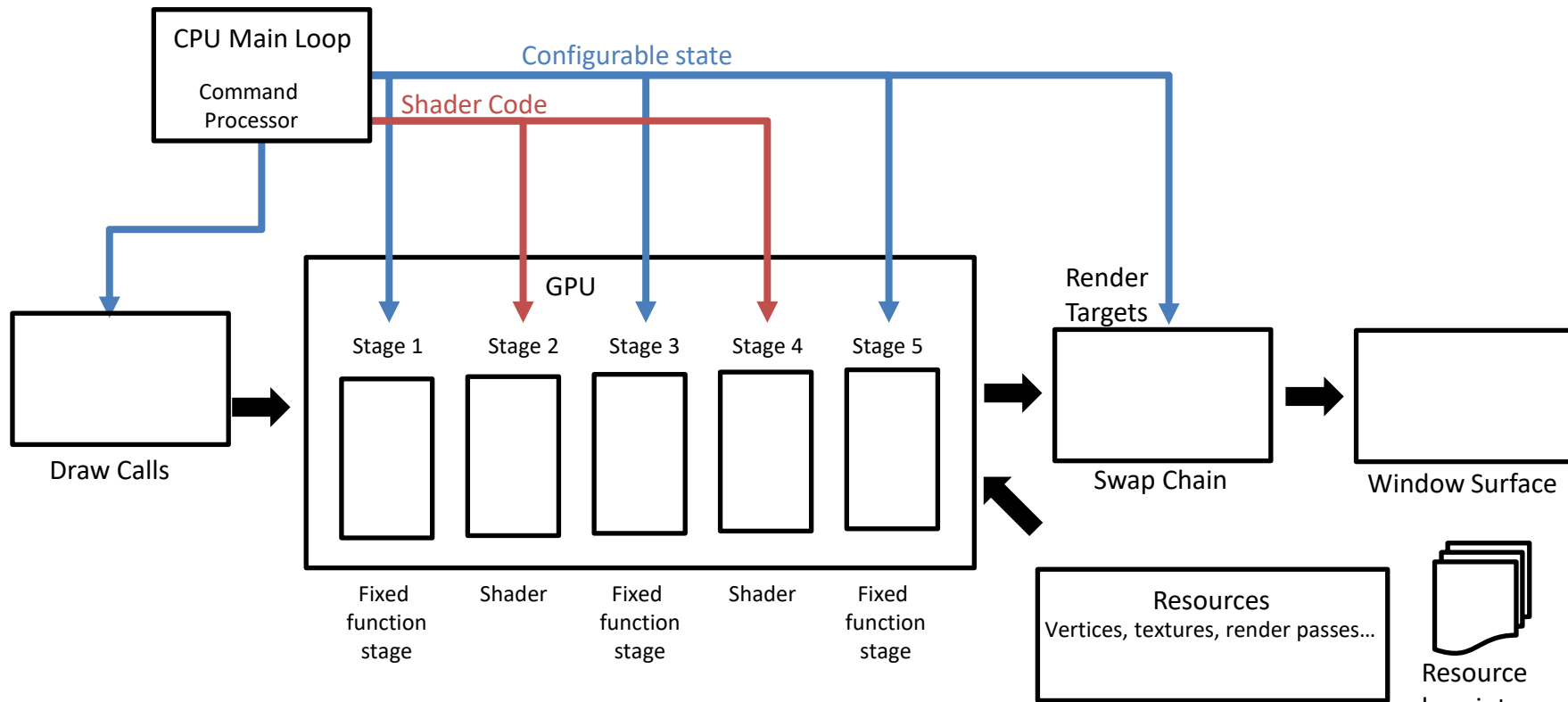- Apple
- ARM (Mali, only design)

# A Story About GPU Programming

- Imagine you want to produce shoes
- CPU, single-threaded
  - The master shoemaker is alone in the shop
  - One task is done after the other
- CPU, multi-threaded
  - The master shoemaker has a few apprentices
  - They talk directly to one another to split the work
- GPU: 10000 worker threads…?

# A Big Shoe Factory

# A Big Graphics API



CPU Main Loop

Command Processor

Configurable state

Shader Code

Draw Calls

GPU

Stage 1 | Stage 2 | Stage 3 | Stage 4 | Stage 5

Fixed function stage | Shader | Fixed function stage | Shader | Fixed function stage

Render Targets

Swap Chain

Window Surface

Resources
Vertices, textures, render passes…
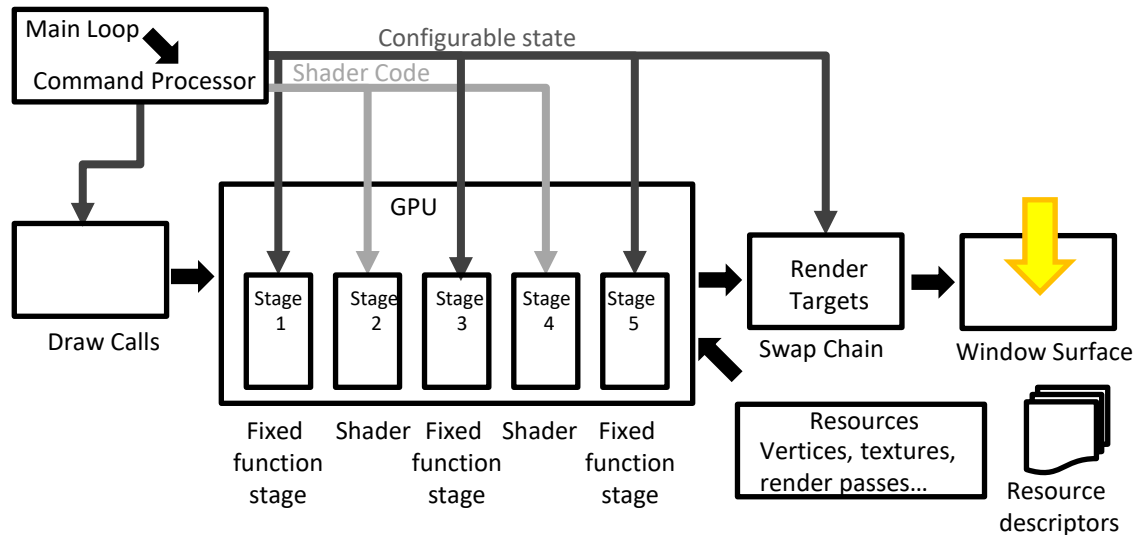
Resource descriptors

# Initialization

- Initialize the graphics system
- Get the "command processor" object
  - Context (OpenGL) or command queue (Vulkan) – see later slides
- Select a render device
  - Which physical GPU
    - If there are multiple
  - Which features are needed
    - Number of viewports, single or double float…
  - Which operation modes are needed
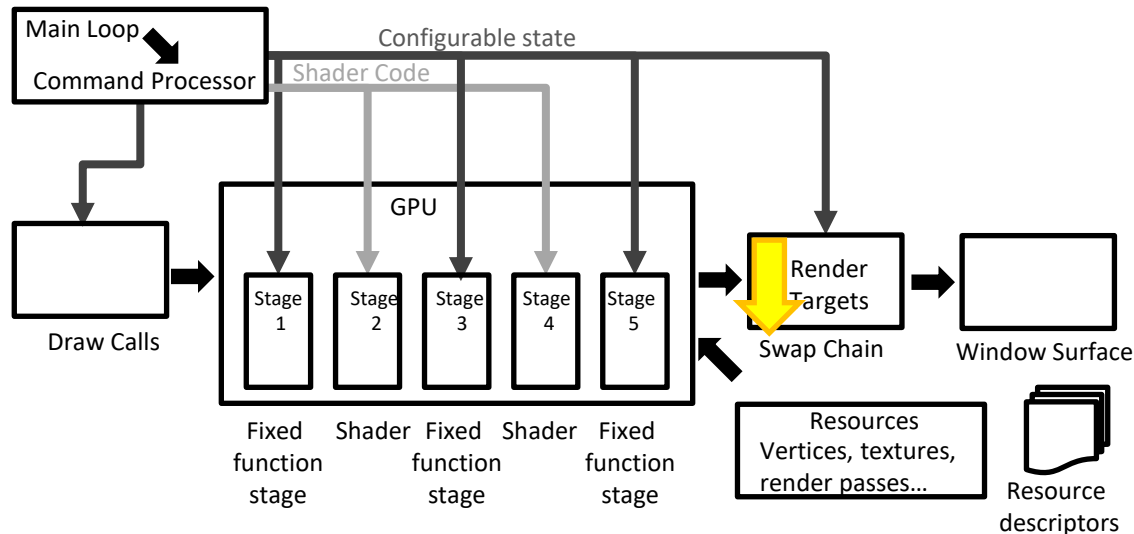    - Graphics, compute, memory transfer

# Surface

- Object connected to window or entire screen
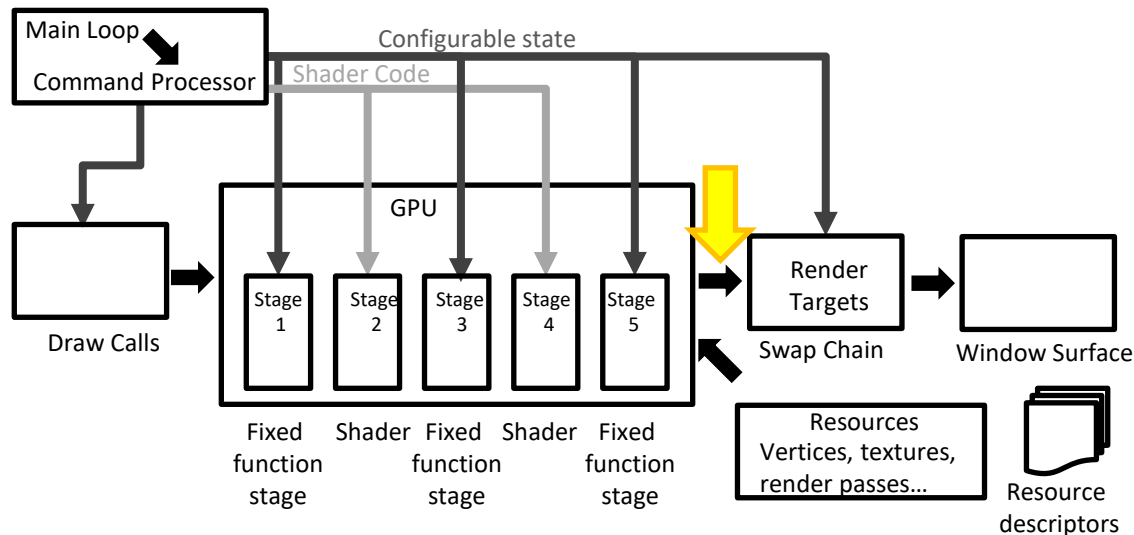- Requires a platform-specific window manager library (e.g., GLFW)

# Swap Chain

- Contains render targets (images) waiting to be shown
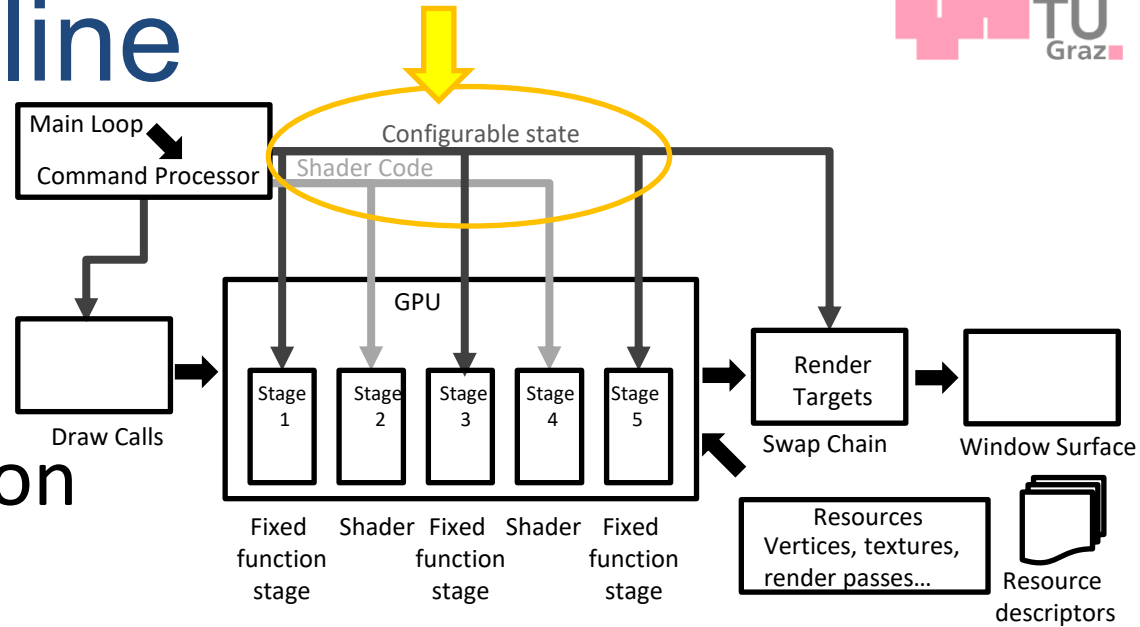- Single, double, triple buffering

# Framebuffer

- Points to render targets for color, depth, stencil
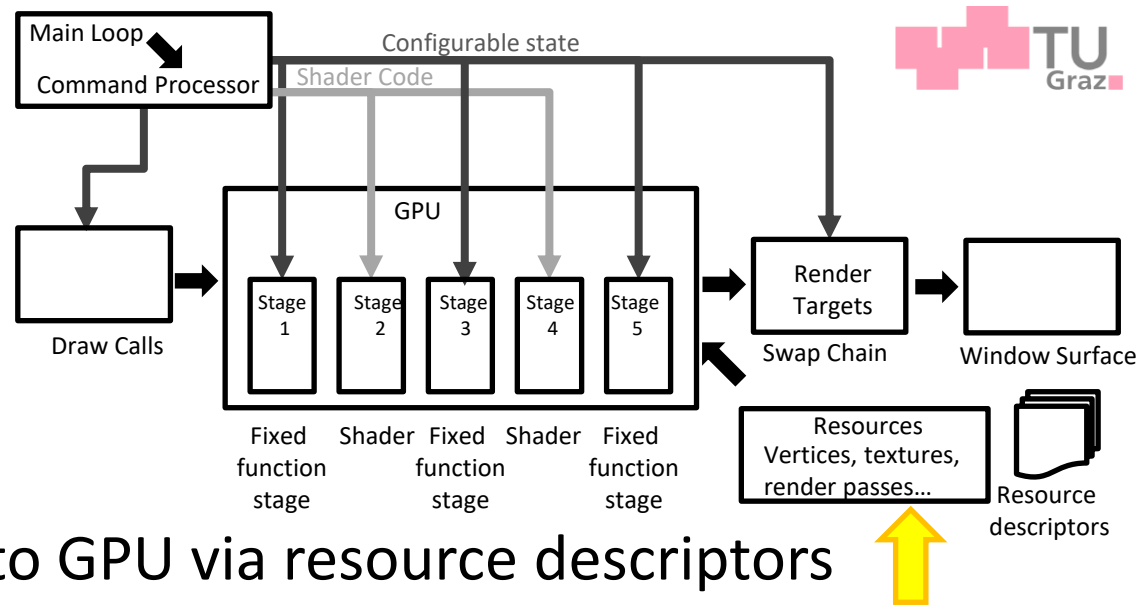
# Graphics Pipeline



- Pipeline stores specific configuration of render state

- Configurable state: viewport size, depth buffer etc.

- Programmable state: shader programs

# Resources



Main Loop
Command Processor
Configurable state
Shader Code
GPU
Draw Calls
Stage 1
Stage 2
Stage 3
Stage 4
Stage 5
Fixed function stage
Shader
Fixed function stage
Shader
Fixed function stage
Render Targets
Swap Chain
Window Surface
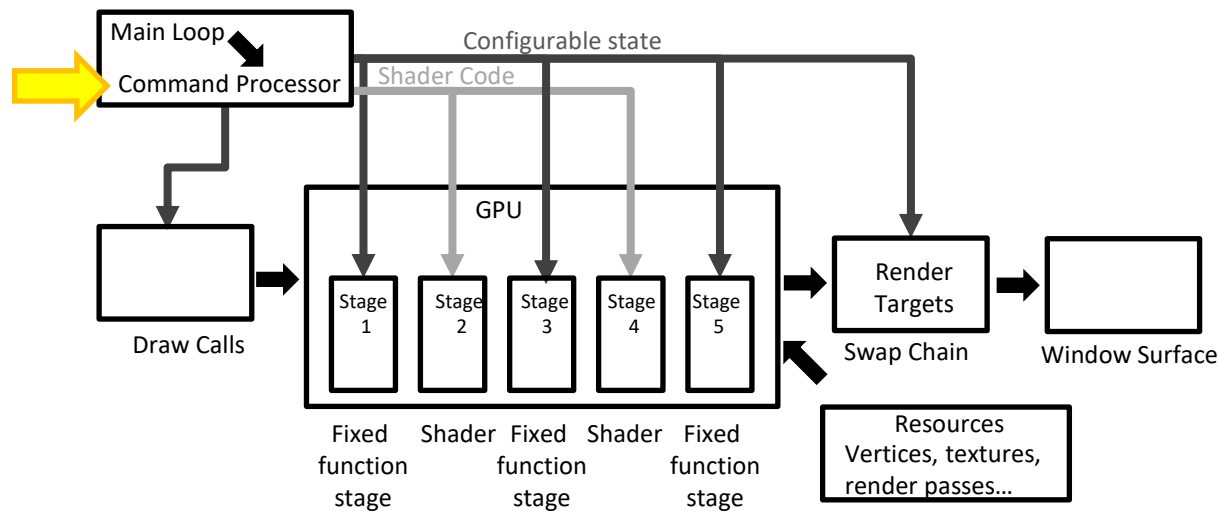Resources Vertices, textures, render passes…
Resource descriptors

- Data stored in GPU local memory

- E.g, vertices, textures, other buffers

- Resource are described to GPU via resource descriptors

- Before use
  - Generate or download resource
  - Specify resource descriptors
  - Bind (=activate) the chosen resource

- Configurable state has "slots" for binding various resource types
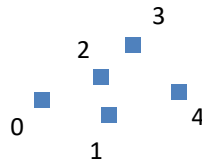
# Command Processor

Example command sequence

- Begin render pass
- Bind pipeline (Vulkan) or pipeline components (OpenGL)
- Bind framebuffer
- Draw vertices (=draw call)
- End render pass

Main Loop

Command Processor

Configurable state

Shader Code

Draw Calls

GPU

Stage 1 — Fixed function stage

Stage 2 — Shader

Stage 3 — Fixed function stage

Stage 4 — Shader

Stage 5 — Fixed function stage

Render Targets

Swap Chain

Window Surface

Resources
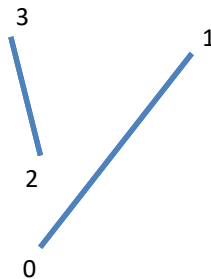Vertices, textures, render passes…

# Draw Calls

- Specify
  - Which primitive type
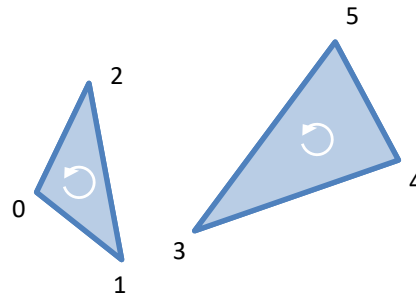  - Which vertex buffer
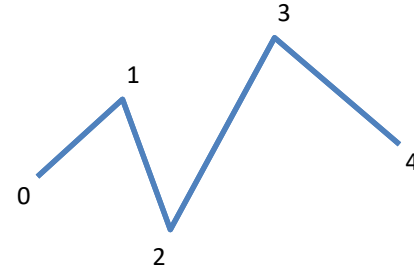  - Start index in buffer
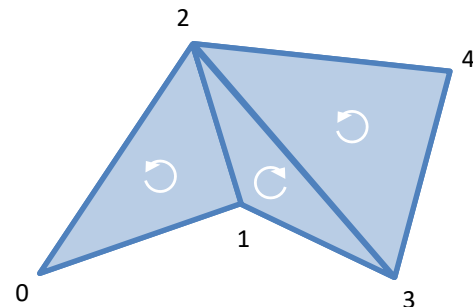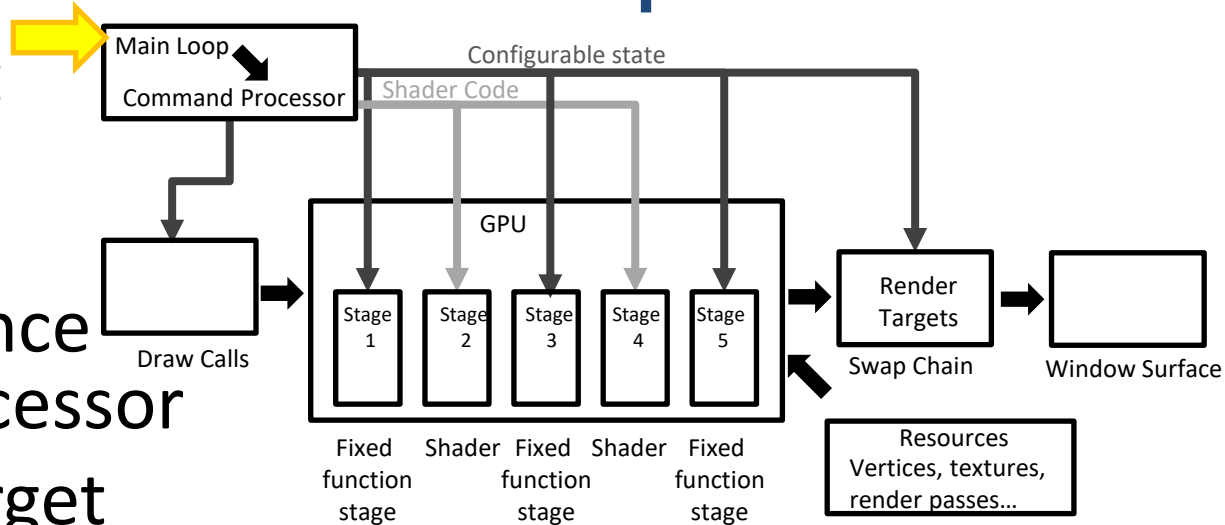  - Stop index in buffer

POINTS

LINES

LINE STRIP

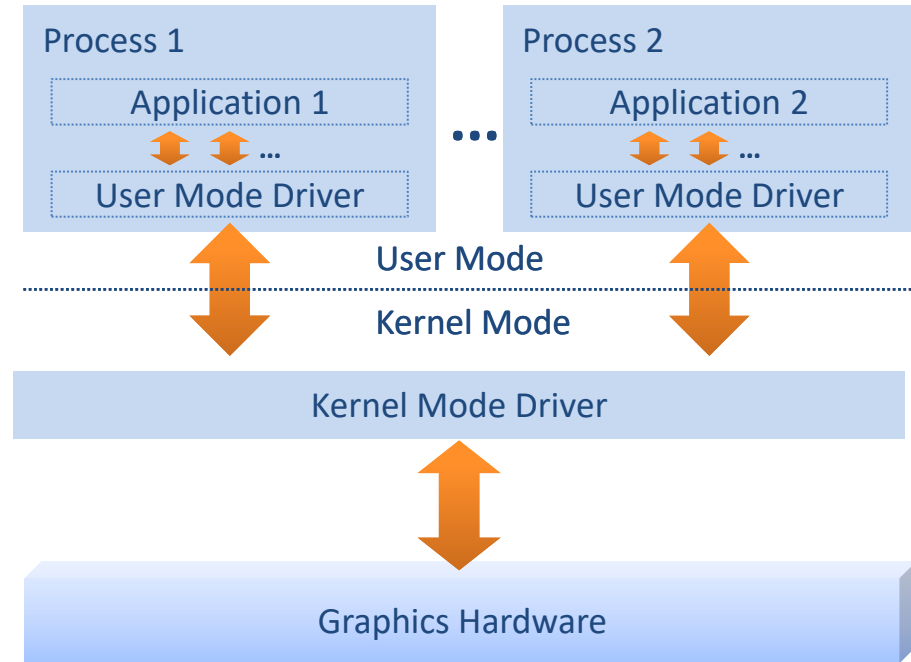TRIANGLES

TRIANGLE STRIP

# CPU Main Loop

- Get render target from swap chain

- Submit chosen command sequence to command processor

- Return render target to swap chain

- Command submission is very different between traditional and modern API

# Graphics Driver Architecture

- User mode graphics driver
  - Minimize number of mode switches
  - Translation of graphics commands to instructions for the hardware
  - Batching, optimization, validation
  - Fine grained memory management
- Kernel mode graphics driver
  - Schedule access to hardware
  - Microkernel pattern, stability
  - Coarse grained memory management
  - Submits command buffers to GPU

Process 1
Application 1
...
User Mode Driver

Process 2
Application 2
...
User Mode Driver

...

User Mode
Kernel Mode

Kernel Mode Driver

Graphics Hardware

# Traditional API

- **Graphics context**: the system object in a traditional API
  - OpenGL, DirectX11 and below
  - Represents a virtual GPU
  - One process can have multiple contexts
  - Multiple contexts can share resources
- **Current context** for a given process
  - One to one mapping
    - Maximum of one current context per thread
    - Current context only assigned to one thread at the same time
  - All OpenGL operations work on current context

# Pipeline State

- Inside a context, one must use commands to configure GPU **pipeline state** before rendering

- Pipeline state consists of

  - Programmable state (shaders)

  - Configurable state: blend, depth, culling, etc.

  - Layout: how to map settings at each stage's shader

# Problems with Traditional API

- Pipeline execution is largely asynchronous
  - CPU sends commands into a "black hole" and
  - CPU does not know exactly when commands are executed
- Configuration of state can only be done incrementally
  - Submitting configuration changes to driver requires immediate validation, conversion, buffering → high cost at runtime
  - Drivers must have per-game optimizations built in
- Pipeline state *not* made explicit in rendering context
  - Switching pipeline configurations is cumbersome
  - Must be done by sequences of state change commands
- Pipeline abstraction is single-threaded on CPU
  - Cannot multi-thread feeding the pipeline
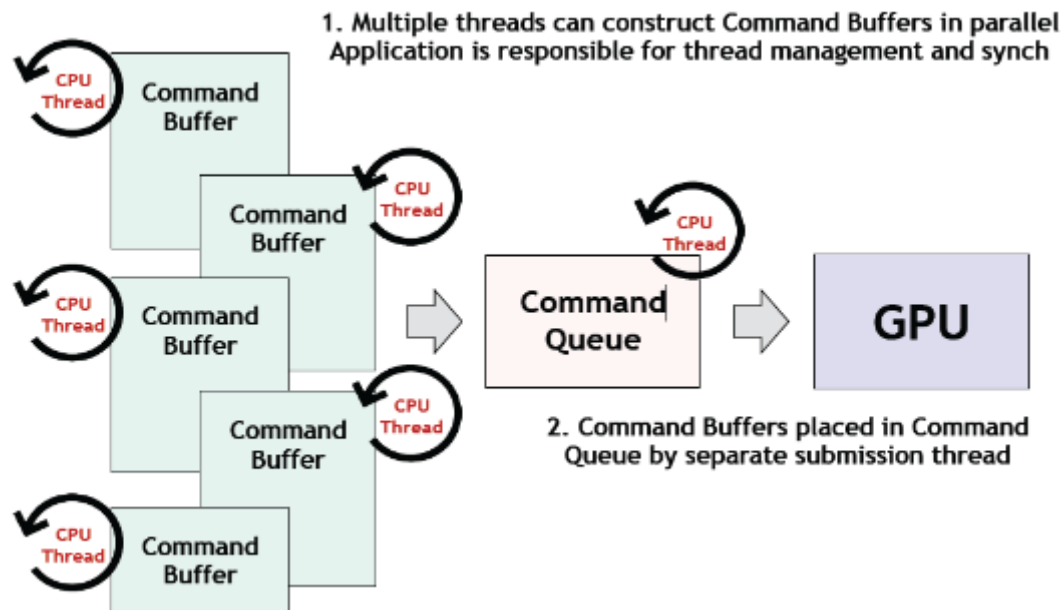  - Having too many draw calls becomes a bottleneck

# Modern API

- DirectX 12, Vulkan (=OpenGL successor)
- Designed for modern GPU types (including mobile)
- Make (CPU driver side of) pipeline more programmable
- CPU multi-threading possible (at your own risk)
- Split rendering context into command buffers and queues
- Make pipeline state (and render passes) explicit
- Low overhead
- Fine-grained control (minimal program 1K lines of code)

# Command Buffers

- Commands collected in command buffers
  - Optimize and validate command buffers during building, not during submission
  - Yields *immutable*, re-useable pipeline state configurations
  - Selected pipeline state variables can be declared *mutable*
- Command buffer submitted to *queue* for execution
- Build *many* command buffers from *many* threads
  - When the buffers are ready, one can submit them all at once
- Each command buffer just *switches* to its favorite pipeline
- Can use *synchronization* primitives across command buffers
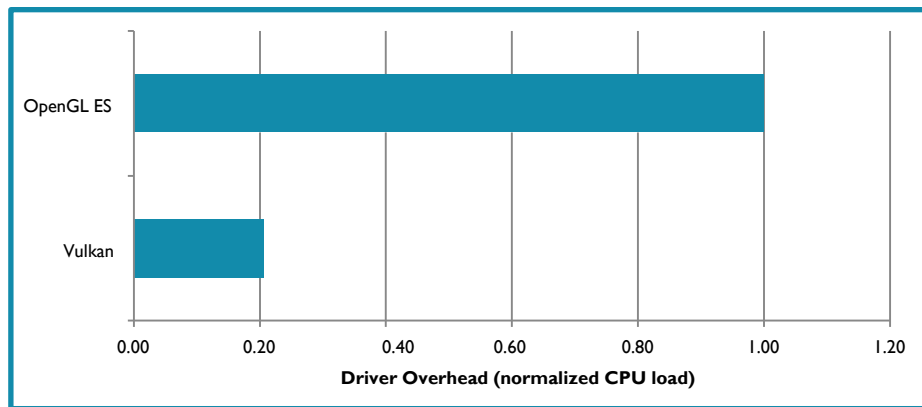  - Event, barrier, semaphore, fence

# Queues

- Queues replace traditional contexts

- Insert command buffer into queue to schedule it

1. Multiple threads can construct Command Buffers in parallel Application is responsible for thread management and synch

CPU Thread — Command Buffer

CPU Thread — Command Buffer

CPU Thread — Command Buffer

CPU Thread — Command Buffer

CPU Thread — Command Buffer

CPU Thread → Command Queue → GPU

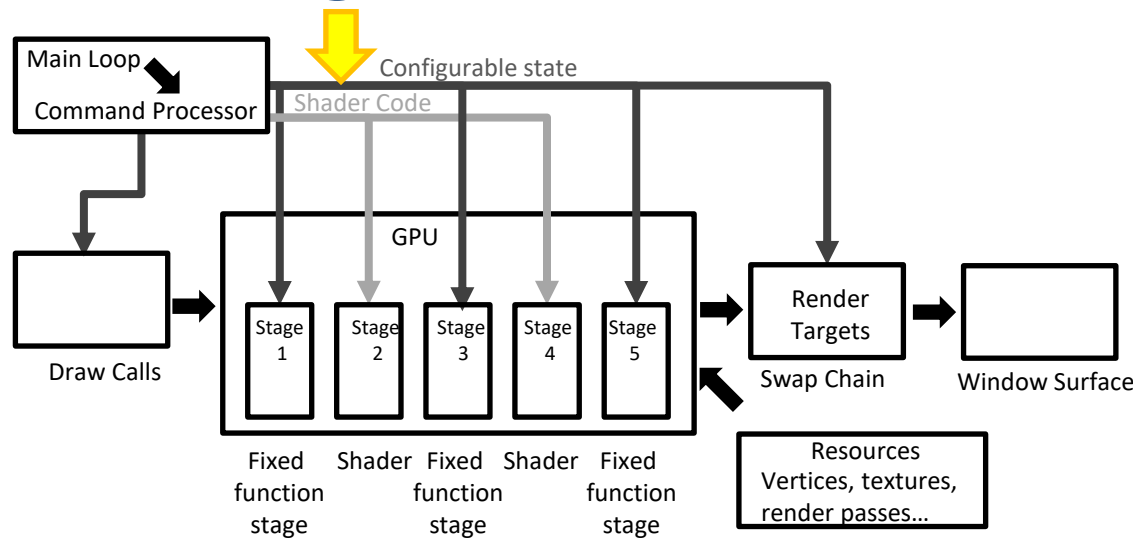2. Command Buffers placed in Command Queue by separate submission thread

# Example: Multi-Threaded Drawing

- ARM Vulkan Benchmark
  - ARM Cortex A-15/7, Mali T-628 MP6
  - 1000 meshes, 3 materials
  - 79% less CPU

Driver Overhead (normalized CPU load)

# Shader Programs

- For the programable parts of the pipeline
- Compiled from shader language
  - HLSL, GLSL
- C-like syntax
- Stream execution model:
  - Shaders are like callbacks, no "main loop" code

Main Loop
Command Processor

Configurable state
Shader Code

GPU

Draw Calls

Stage 1 | Stage 2 | Stage 3 | Stage 4 | Stage 5

Fixed function stage | Shader | Fixed function stage | Shader | Fixed function stage

Render Targets

Swap Chain

Window Surface

Resources
Vertices, textures, render passes…

# Language Elements

- Skalar data types: float, int, bool...
- Vector/matrix data types: vec2, vec3, vec4, ivec3, mat4...
- Struct, arrays (static size)
- Texture sampler: sample2D, sampler3D, ...
- Control flow: if, else, while, for
- Function calls (no recursion)
- Swizzle operators:   color1.rgb = color2.bga
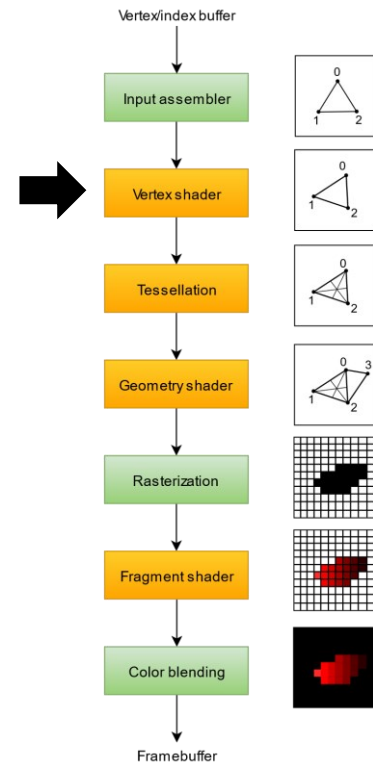- Mask operator: pos1.z = pos2.x + pos2.y

# Build-In Functions

- Trigonometric
  - sin(), cos(), radians(), …
- Logarithm, exponentiation
  - log(), sqrt(), …
- Other
  - min(), max(), mod(), floor(), abs(), …
- Geometric
  - distance(), normalize(), dot(), length(), …
- Special functions
  - Linear interpolation
  - Reflection vector
  - Refraction vector

# Shader Compilation

- Compilation separated into front-end/back-end
- Front-end: HLSL, GLSL, OpenCL, etc.
- Back-end = Bytecode
  - Microsoft HLSL uses proprietary bytecoe
  - Vulkan uses SPIR-V (standard portable intermediate representation)
- Applications ship with bytecode, not shader source
- Just-in-time compilation of bytecode
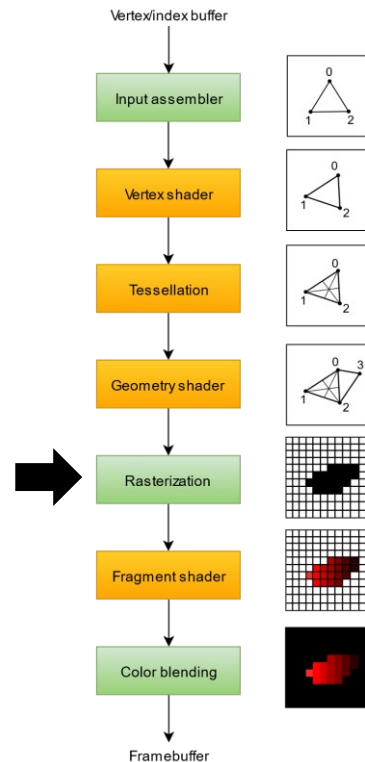  - To hardware target platform (NVIDIA, AMD, …)

# Vertex Shader

- Processes each vertex

- Input: vertex attributes

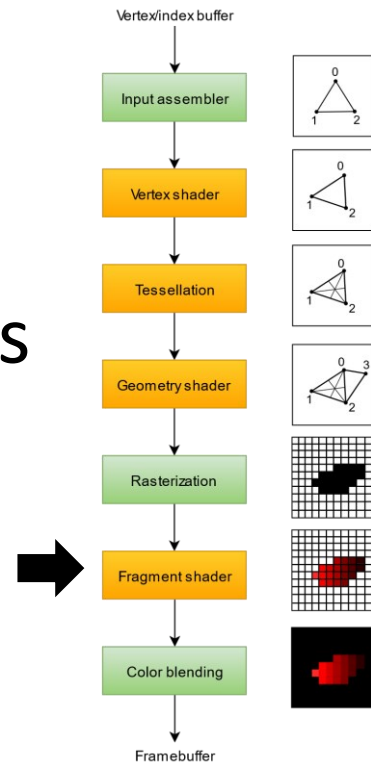- Output: vertex attributes

- Mandatory output: gl_Position

# Rasterizer

- Generates fragments covering a primitive

- Fixed-function unit

- Input: primitives, vertex attributes

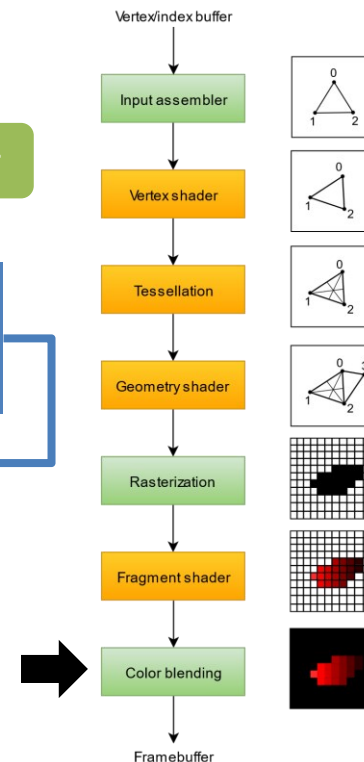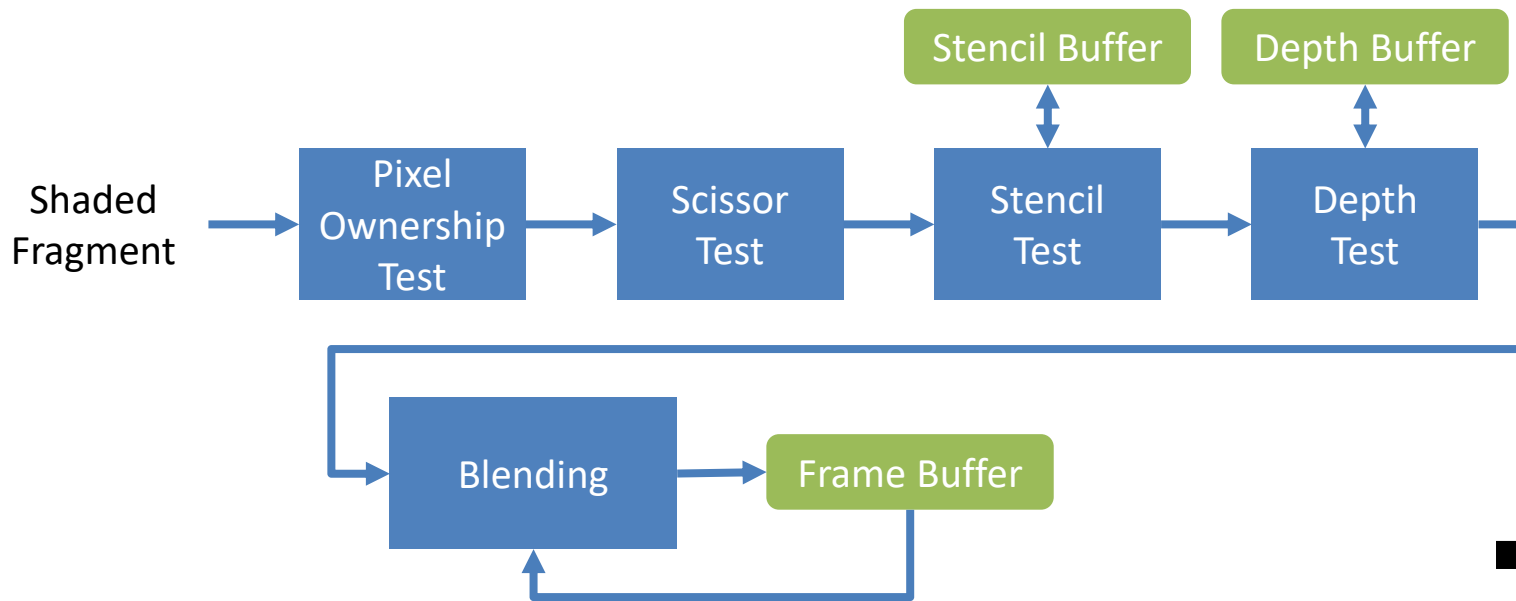- Output: fragments with interpolated vertex attributes

# Fragment Shader



- Processes each fragment

- Input: interpolated vertex attributes

- Output: fragment color

# Fragment Merging



Shaded Fragment → Pixel Ownership Test → Scissor Test → Stencil Test → Depth Test

Stencil Buffer ↕ Stencil Test

Depth Buffer ↕ Depth Test

Blending → Frame Buffer

# Anatomy of a GLSL Shader

```glsl
1  #version 330

2  // uniform inputs are constant for all shader invocations
3  uniform vec4 some_uniform;

4  // inputs are varying with each shader invocation
5  layout(location = 0) in vec3 some_input;
6  layout(location = 1) in vec4 another_input;

7  // outputs
8  out vec4 some_output;

9  void main()
10 {
11     //…
12 }
```

# Built-In Variables

- Interface to fixed-function parts of pipeline
  - E. g. vertex shader:
    - `in int gl_VertexID;`
    - `out vec4 gl_Position;`
  - E. g. fragment shader:
    - `in vec4 gl_FragCoord;`
    - `out float gl_FragDepth;`

# Example: Vertex Shader

```
1  #version 330

2  uniform mat4 mvp_matrix; // model-view-projection matrix

3  layout(location = 0) in vec3 vertex_position;
4  layout(location = 1) in vec4 vertex_color;

5  out vec4 color;

6  void main()
7  {
8      gl_Position = mvp_matrix * vec4(vertex_position, 1.0f);
       color = vertex_color;
9  }
```

# Example: Fragment Shader

```
1  #version 330

2  layout(location = 0) in vec4 color; // interpolated color

3  out vec4 fragment_color;

4  void main()
5  {
6      fragment_color = color;
7  }
```

# Thank You!

# Questions ?