

Vulkan.
Introduction

Oregon State University
Mike Bailey
mb@cs.oregonstate.edu

CC BY-NC-ND
This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

Oregon State University Computer Graphics info.2020 mb - September 17, 2018

Acknowledgements



First of all, thanks to the inaugural class of 19 students who braved new, unrefined, and just-in-time course materials to take the first Vulkan class at Oregon State University – Winter Quarter, 2018. Thanks for your courage and patience!

Second, thanks to NVIDIA! The GeForce 1080Ti cards are what made this course possible.

Third, thanks to Kathleen Mattison and the Khronos Group for the great laminated Vulkan Quick Reference Cards! (Look at those happy faces in the photo holding them.)

Oregon State University
NVIDIA
KHRONOS GROUP

Oregon State University Computer Graphics mb - September 17, 2018

What Prompted the Move to Vulkan?

1. Performance
2. Performance
3. Performance

Vulkan is better at keeping the GPU busy than OpenGL is. OpenGL drivers need to do a lot of CPU work before handing work off to the GPU. Vulkan lets you get more power from the GPU card you already have.

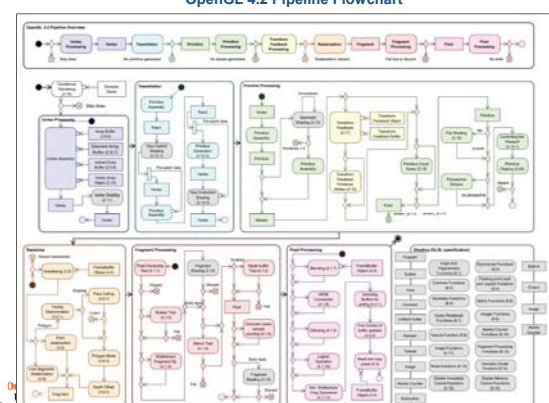
This is especially important if you can hide the complexity of Vulkan from your customer base and just let them see the improved performance. Thus, Vulkan has had a lot of support and interest from game engine developers, 3rd party software vendors, etc.

Oregon State University Computer Graphics

As an aside, the Vulkan development effort was originally called **glNext**, which created the false impression that this was a replacement for OpenGL. It's not.

mb - September 17, 2018

OpenGL 4.2 Pipeline Flowchart



Oregon State University Computer Graphics mb - September 17, 2018

Why is it so important to keep the GPU Busy?

Nvidia Titan V Specs vs. Titan Xp, 1080 Ti					
	Titan V	Titan V100	Titan P100	GTX 1080 Ti	GTX 1080
GPU	GP100	GP100	GP100 Cut Down Pascal	GP102 Pascal	GP104-400 Pascal
Transistor Count	21.1B	21.1B	15.3B	12B	7.2B
Full Process	12nm FFN	12nm FFN	16nm FinFET	16nm FinFET	16nm FinFET
CUDA Cores / Tensor Cores	5120 / 640	5120 / 640	3584 / 0	3584 / 0	2560 / 0
TB/s	320	324	224	224	160
TP/s	?	?	96 (T)	88	64
Core Clock	1500MHz	1500MHz	1500MHz	1600MHz	1600MHz
Boost Clock	1620MHz	1730MHz	1620MHz	1600MHz	1730MHz
FP32 TLOps	16371LOps	14371LOps	16371LOps	~14371LOps	9171LOps
Memory Type	HBM2	HBM2	HBM2	GDDR5X	GDDR5X
Memory Capacity	12GB	16GB	16GB	11GB	6GB
Memory Clock	1.70Gps HBM2	1.75Gps HBM2	?	11Gps	10Gps GDDR5X
HBM22 Bandwidth	312TB/s	406TB/s	406TB/s	352TB/s	256TB/s
Memory Bandwidth	652GB/s	900GB/s	?	~486GB/s	303.3GB/s
Total Power Budget (TDP)	250W	250W	300W	250W	180W
Power Connectors	1x 8-pin 1x 6-pin	?	?	1x 8-pin 1x 6-pin	1x 8-pin
Release Date	12/07/2017	02/15/2017	1B3	1B3	02/20/16
Release Price	\$3000	\$10000	-	\$700	MSRP: \$600 New \$300

Oregon State University Computer Graphics

The Nvidia Titan V graphics card is not targeted at gamers, but rather at scientific and machine-learning workloads. That does not, however, mean that the card is incapable of gaming, nor does it mean that we can't extrapolate future key performance metrics for Vulkan. The Titan V is a derivative of the earlier released GV100 (GPU) part of the Tesla architecture card series. The key difference is that the Titan V ships at \$3000, whereas the Tesla V100 was available as part of a \$10,000 developer kit. The Tesla V100 still offers greater memory capacity by a factor of 2 (~16GB HBM2 versus 12GB HBM2) and has a wider memory interface, but other core features remain matched or nearly matched. Core clock, for one, is 1500 MHz versus 1600 MHz, with 640 tensor cores (used for TensorFlow deep-learning workloads) on each GPU.


mb - September 17, 2018

Who was the original Vulcan?

From Wikipedia:

"Vulcan is the god of fire including the fire of volcanoes, metalworking, and the forge in ancient Roman religion and myth. Vulcan is often depicted with a blacksmith's hammer. The **Vulcanalia** was the annual festival held August 23 in his honor. His Greek counterpart is Hephaestus, the god of fire and smithery. In Etruscan religion, he is identified with Sethlans. Vulcan belongs to the most ancient stage of Roman religion: Varro, the ancient Roman scholar and writer, citing the *Annales Maximi*, records that king Titus Tatius dedicated altars to a series of deities among which Vulcan is mentioned."

[https://en.wikipedia.org/wiki/Vulcan_\(mythology\)](https://en.wikipedia.org/wiki/Vulcan_(mythology))



Oregon State University Computer Graphics mb - September 17, 2018

Why Name it after the God of the Forge?

Oregon State University Computer Graphics
mp - September 17, 2018

Who is the Khronos Group?

The Khronos Group, Inc. is a non-profit member-funded industry consortium, focused on the creation of open standard, royalty-free application programming interfaces (APIs) for authoring and accelerated playback of dynamic media on a wide variety of platforms and devices. Khronos members may contribute to the development of Khronos API specifications, vote at various stages before public deployment, and accelerate delivery of their platforms and applications through early access to specification drafts and conformance tests.

Oregon State University Computer Graphics
mp - September 17, 2018

Playing "Where's Waldo" with Khronos Membership

Oregon State University Computer Graphics
mp - September 17, 2018

Who's Been Specifically Working on Vulkan?

Oregon State University Computer Graphics
mp - September 17, 2018

Vulkan

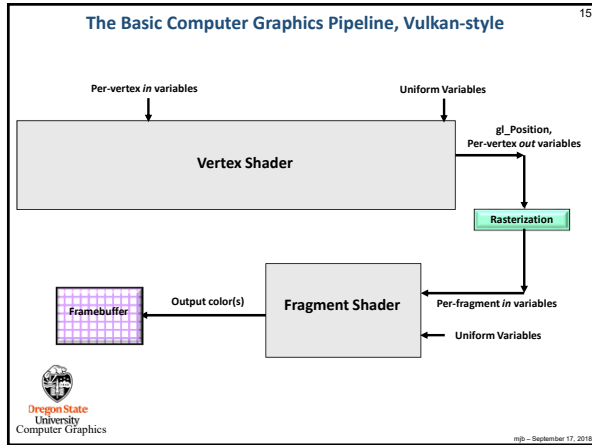
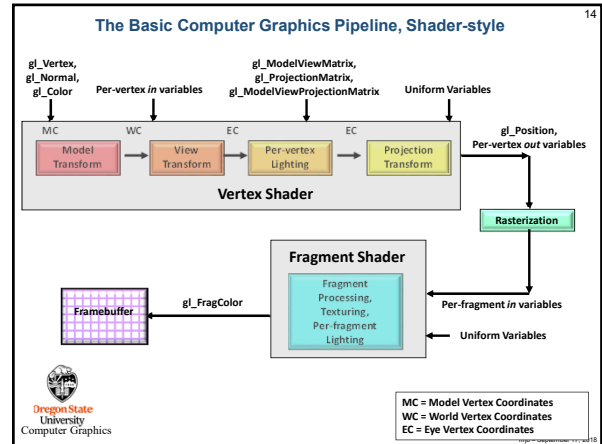
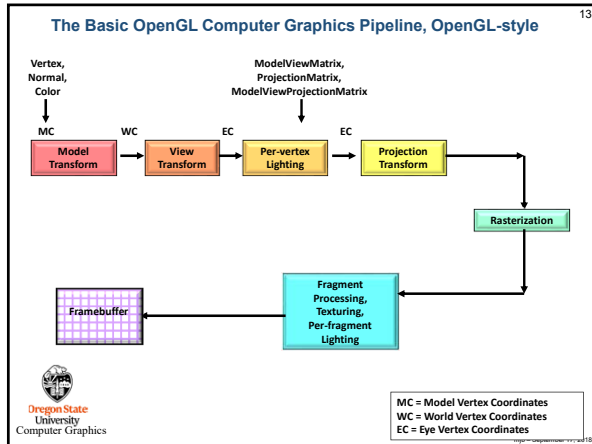
- Largely derived from AMD's Mantle API
- Also heavily influenced by Apple's Metal API and Microsoft's DirectX 12
- Goal: much less driver complexity and overhead than OpenGL has
- Goal: much less user hand-holding – Vulkan can crash
- Goal: higher single-threaded performance than OpenGL can deliver
- Goal: able to do multithreaded graphics
- Goal: able to handle tiled rendering

Oregon State University Computer Graphics
mp - September 17, 2018

Vulkan Differences from OpenGL

- More low-level information must be provided (by you!) in the application, rather than the driver
- Screen coordinate system is Y-down
- No "current state", at least not one maintained by the driver
- All of the things that we have talked about being *deprecated* in OpenGL are *really deprecated* in Vulkan: built-in pipeline transformations, begin-end, fixed-function, etc.
- You must manage your own transformations.
- All transformation, color, texture functionality must be done in shaders.
- Shaders are pre-"half-compiled" outside of your application. The compilation process is then finished during the pipeline-building process.

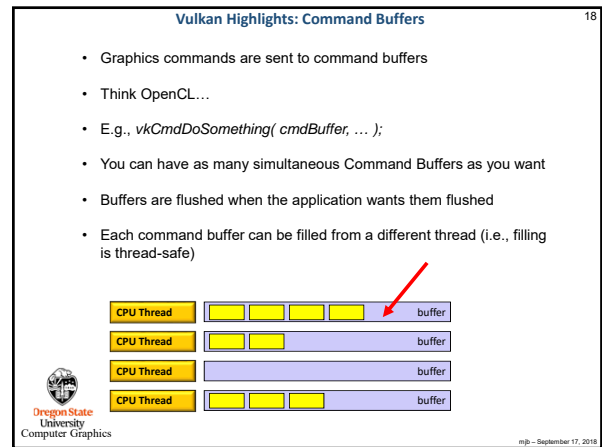
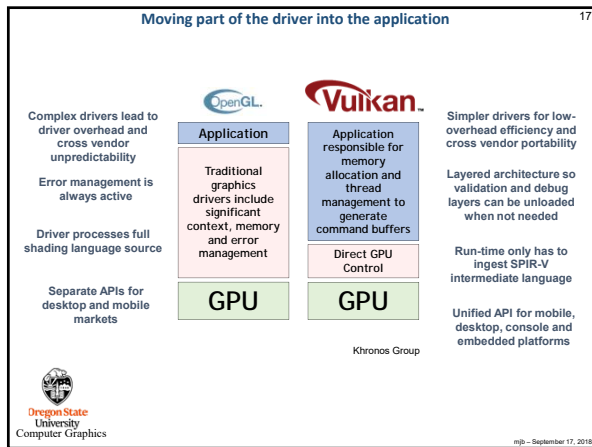
Oregon State University Computer Graphics
mp - September 17, 2018



A Complete API Redesign

OpenGL	Vulkan
Originally architected for graphics workstations with direct renderers and split memory	Matches architecture of modern platforms including mobile platforms with unified memory, tiled rendering
Driver does lots of work: state validation, dependency tracking, error checking. Limits and randomizes performance	Explicit API – the application has direct, predictable control over the operation of the GPU
Threading model doesn't enable generation of graphics commands in parallel to command execution	Multi-core friendly with multiple command buffers that can be created in parallel
Syntax evolved over twenty years – complex API choices can obscure optimal performance path	Removing legacy requirements simplifies API design, reduces specification size and enables clear usage guidance
Shader language compiler built into driver. Only GLSL supported. Have to ship shader source	SPIR-V as compiler target simplifies driver and enables front-end language flexibility and reliability
Despite conformance testing, developers must often handle implementation variability between vendors	Simpler API, common language front-ends, more rigorous testing increase cross vendor functional/performance portability


Khronos Group



Vulkan Highlights: Pipelines

19

- In OpenGL, your "pipeline state" is whatever your current graphics attributes are: color, transformations, textures, shaders, etc.
- Changing the state on-the-fly one item at-a-time is very expensive
- Vulkan forces you to set all your state at once into a "pipeline state object" (PSO) and then invoke the entire PSO whenever you want to use that state combination
- Think of the pipeline state as being immutable.
- Potentially, you could have thousands of these pre-prepared states
- This is a good time to talk about how game companies view Vulkan...



mjb - September 17, 2018

Querying the Number of Things

20


```
uint32_t count;
result = vkEnumeratePhysicalDevices( Instance, OUT &count, OUT (VkPhysicalDevice *) nullptr );

VkPhysicalDevice * physicalDevices = new VkPhysicalDevice[ count ];
result = vkEnumeratePhysicalDevices( Instance, OUT &count, OUT physicalDevices );
```

This way of querying information is a recurring OpenGL and Vulkan pattern (get used to it):

```
result = vkEnumeratePhysicalDevices( Instance, &count, nullptr );
result = vkEnumeratePhysicalDevices( Instance, &count, physicalDevices );
```

How many total there are Where to put them



mjb - September 17, 2018

Vulkan Code has a Distinct "Style"

21


```
VkBufferCreateInfo vbci;
vbci.sType = VK_STRUCTURE_TYPE_BUFFER_CREATE_INFO;
vbci.pNext = nullptr;
vbci.flags = 0;
vbci.size = << buffer size in bytes >>
vbci.usage = VK_USAGE_UNIFORM_BUFFER_BIT;
vbci.sharingMode = VK_SHARING_MODE_EXCLUSIVE;
vbci.queueFamilyIndexCount = 0;
vbci.pQueueFamilyIndices = nullptr;

VK_RESULT result = vkCreateBuffer ( LogicalDevice, IN &vbci, PALLOCATOR, OUT &Buffer );

VkMemoryRequirements vmr;
result = vkGetBufferMemoryRequirements( LogicalDevice, Buffer, OUT &vmr ); // fills vmr

VkMemoryAllocateInfo vmai;
vmai.sType = VK_STRUCTURE_TYPE_MEMORY_ALLOCATE_INFO;
vmai.pNext = nullptr;
vmai.flags = 0;
vmai.allocationSize = vmr.size;
vmai.memoryTypeIndex = 0;


result = vkAllocateMemory( LogicalDevice, IN &vmai, PALLOCATOR, &MatrixBufferMemoryHandle );
result = vkBindBufferMemory( LogicalDevice, Buffer, MatrixBufferMemoryHandle, 0 );
```



mjb - September 17, 2018

Vulkan Quick Reference Card

22

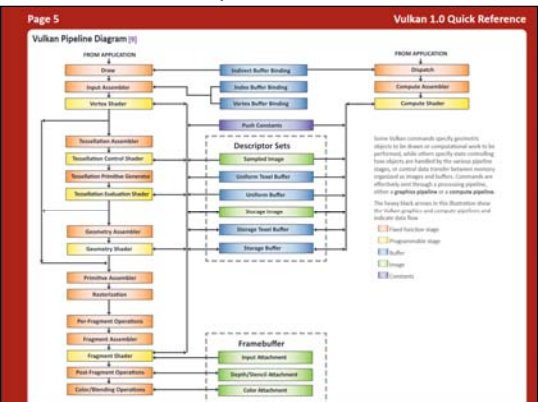


<https://www.khronos.org/files/vulkan1-reference-guide.pdf>

mjb - September 17, 2018

Vulkan Quick Reference Card

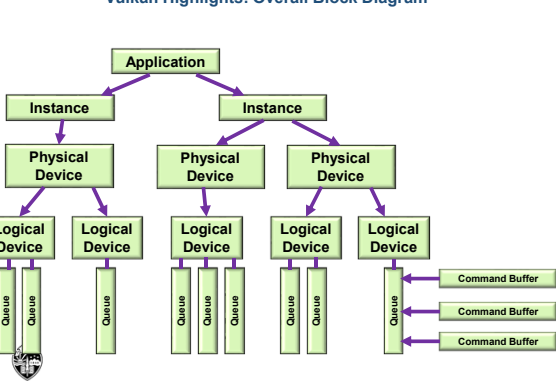

23



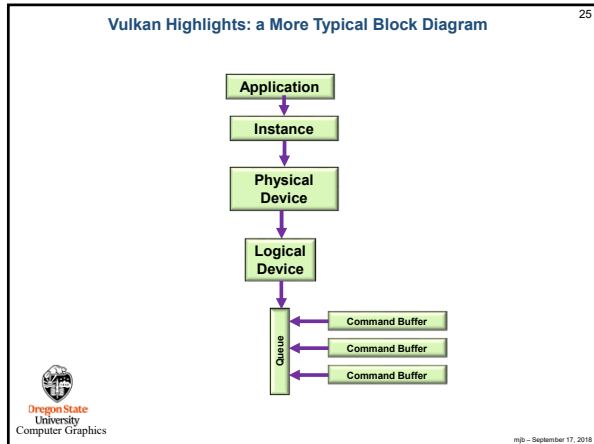
mjb - September 17, 2018

Vulkan Highlights: Overall Block Diagram

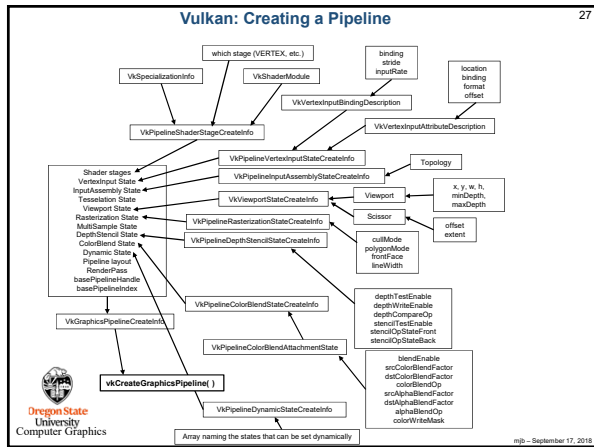
24

mjb - September 17, 2018



- ### Steps in Creating Graphics using Vulkan
1. Create the Instance
 2. Setup the Debug Callbacks
 3. Create the Surface
 4. List the Physical Devices
 5. Pick the right Physical Device
 6. Create the Logical Device
 7. Create the Uniform Variable Buffers
 8. Create the Vertex Data Buffers
 9. Create the texture sampler
 10. Create the texture images
 11. Create the Swap Chain
 12. Create the Depth and Stencil Images
 13. Create the RenderPass
 14. Create the Framebuffer(s)
 15. Create the Descriptor Set Pool
 16. Create the Command Buffer Pool
 17. Create the Command Buffer(s)
 18. Read the shaders
 19. Create the Descriptor Set Layouts
 20. Create and populate the Descriptor Sets
 21. Create the Graphics Pipeline(s)
 22. Update-Render-Update-Render- ...
- Oregon State University
Computer Graphics
- 26



Vulkan GPU Memory

- Your application allocates GPU memory for the objects it needs
- You map GPU memory to the CPU address space for access
- Your application is responsible for making sure what you put into that memory is actually in the right format, is the right size, has the right alignment, etc.

From the OpenGL Shader Storage Buffer notes:

```

glGenBuffers( 1, &posSsbo );
glBindBuffer( GL_SHADER_STORAGE_BUFFER, posSsbo );
glBufferData( GL_SHADER_STORAGE_BUFFER, NUM_PARTICLES * sizeof(struct pos), NULL, GL_STATIC_DRAW );
glInti bufMask = GL_MAP_WRITE_BIT | GL_MAP_INVALIDATE_BUFFER_BIT ; // the invalidate makes a big difference when re-writing
struct pos *points = (struct pos *) glMapBufferRange( GL_SHADER_STORAGE_BUFFER, 0, NUM_PARTICLES * sizeof(struct pos), bufMask );
    
```

Oregon State University
Computer Graphics

28


- ### Vulkan Render Passes
- Drawing is done inside a render pass
 - Each render pass contains what framebuffer attachments to use
 - Each render pass is told what to do when it begins and ends
 - Multiple render passes can be merged
- Oregon State University
Computer Graphics
- 29

- ### Vulkan Compute Shaders
- Compute pipelines are allowed, but they are treated as something special (just like OpenGL does)
 - Compute passes are launched through dispatches
 - Compute command buffers can be run asynchronously
- Oregon State University
Computer Graphics
- 30

Vulkan Synchronization

31

- Vulkan tries to run "flat out"
- Therefore, synchronization is the responsibility of the application
- Events can be set, polled, and waited for (much like OpenCL)
- Vulkan does not ever lock – that's the application's job
- Threads can concurrently read from the same object
- Threads can concurrently write to different objects

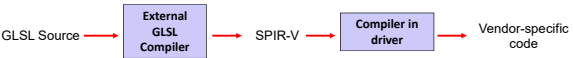


m@ - September 17, 2018

Vulkan Shaders

32

- GLSL is the same as before ... almost
- For places it's not, an implied `#define VULKAN 100` is automatically supplied by the compiler
- You pre-compile your shaders with an external compiler
- Your shaders get turned into an intermediate form known as SPIR-V (Standard Portable Intermediate Representation for Vulkan)
- SPIR-V gets turned into fully-compiled code at runtime
- The SPIR-V spec has been public for months –new shader languages are surely being developed
- OpenCL and OpenGL will be moving to SPIR-V as well

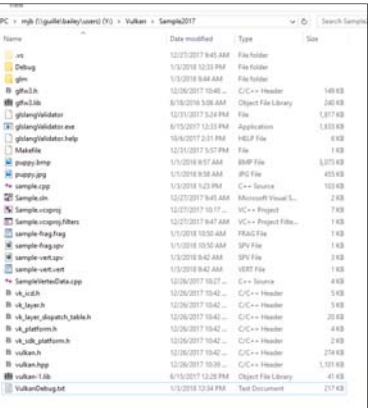


Advantages:


1. Software vendors don't need to ship their shader source
2. Software can launch faster because half of the compilation has already taken place
3. This guarantees a common front-end syntax
4. This allows for other language front-ends

Your Sample2017.zip File Contains This

33



Name	Date modified	Type	Size
ms	12/27/2017 9:45 AM	File folder	
Debug	1/3/2018 12:33 PM	File folder	
glm	1/3/2018 9:48 AM	File folder	
glsl.h	12/26/2017 10:40...	C/C++ Header	149 KB
glsl.h.in	6/76/2016 10:48 AM	Object File Library	240 KB
glslangValidator	12/21/2017 3:38 PM	File	1,917 KB
glslangValidator.exe	6/75/2017 12:33 PM	Application	1,813 KB
glslangValidator.help	10/4/2017 3:31 PM	HELP File	4 KB
Makefile	12/31/2017 3:27 PM	File	1 KB
proprietary	1/3/2018 9:57 AM	BIFF File	3,374 KB
proprietary.jpg	1/3/2018 9:58 AM	JPG File	415 KB
sample.cpp	1/3/2018 1:23 PM	C++ Source	103 KB
sample.h	12/27/2017 9:45 AM	Microsoft Visual S...	2 KB
sample.cppproj	12/27/2017 9:47...	VC++ Project	7 KB
sample.cppproj.filters	12/27/2017 9:47 AM	VC++ Project File...	1 KB
sample.frag.frag	1/3/2018 10:50 AM	FRAG File	1 KB
sample.frag.spr	1/3/2018 10:50 AM	SPR File	1 KB
sample.vert.spr	1/3/2018 9:42 AM	SPR File	1 KB
sample.vert.vert	1/3/2018 9:42 AM	VERT File	1 KB
SampleMainData.cpp	12/26/2017 10:17...	C++ Source	4 KB
vk.h	12/26/2017 10:42...	C/C++ Header	5 KB
vk_base.h	12/26/2017 10:42...	C/C++ Header	5 KB
vk_layer_dispatch_table.h	12/26/2017 10:42...	C/C++ Header	20 KB
vk_layer.h	12/26/2017 10:42...	C/C++ Header	4 KB
vk_layer_dispatch_table.h	12/26/2017 10:42...	C/C++ Header	2 KB
vulkan.h	12/26/2017 10:42...	C/C++ Header	274 KB
vulkan.hpp	12/26/2017 10:39...	C/C++ Header	1,101 KB
vulkan.lib	6/75/2017 12:28 PM	Object File Library	41 KB
vulkanDebug.txt	1/3/2018 12:38 PM	Text Document	217 KB



m@ - September 17, 2018