

# Embedded Outreach

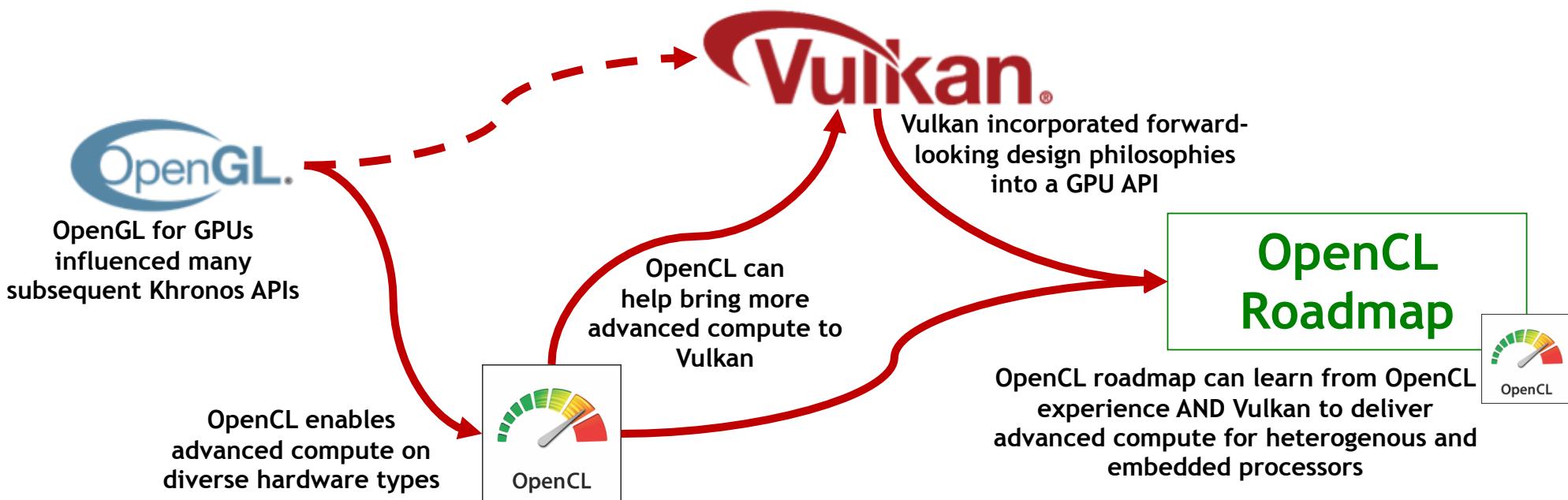
## Focusing OpenCL for the Embedded Industry

### January 2018

Neil Trevett | Khronos President  
NVIDIA | VP Developer Ecosystem  
[ntrevett@nvidia.com](mailto:ntrevett@nvidia.com) | [@neilt3d](https://twitter.com/neilt3d)  
[www.khronos.org](http://www.khronos.org)

# Agenda for this Session

- Introduction to Khronos and the need for open standards
- Update on OpenCL (and SPIR-V) - the standards for heterogenous computing
- Update on Vulkan - a new explicit-style API
- Lessons from Vulkan and OpenCL that can guide the OpenCL roadmap
- Directions for OpenCL roadmap - and increasing relevance to the embedded industry



# What is An Open Standard?



An INTEROPERABILITY STANDARD  
enables two interoperating entities to  
COMMUNICATE

Software <-> Hardware  
Sensor <-> Processor  
Device <-> Host



## Bad Standards

- Overprescribes implementation details
- Forces everyone to implement a lowest common denominator
  - Stifles innovation
- > Commoditization



## Good Standards

- Prescribes only interoperability
  - Enables implementation diversity
  - Encourages innovation
- > Differentiation



## A truly OPEN standard

Is not controlled by a single company - but by the whole industry  
Is freely available to use by any company without royalty payments

# Why Do We Need Open Standards?

An advanced product such as mobile uses *100s* of standards

Wireless: LTE, WiFi, Bluetooth

Physical: HDMI, USB, CSI

Software APIs: HTML5, OpenGL ES

Standards drive mobile market growth by expanding device capabilities



## Standard Defining Organizations



**Standards Grow Markets...**  
By reducing consumer confusion and increasing usability

**... Reduce Costs ...**

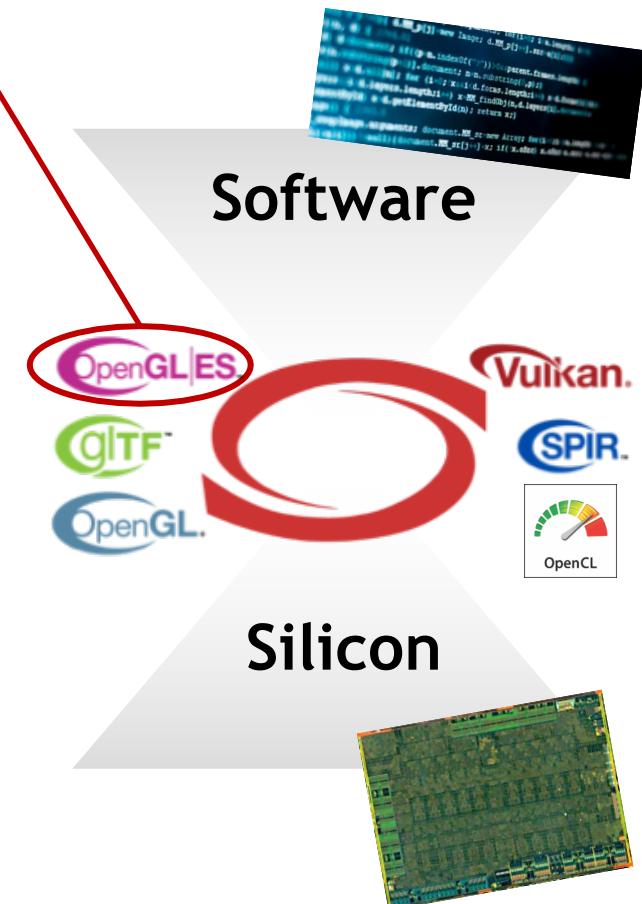
By sharing development between many companies and driving volume

**... Accelerate Time to Market**  
With well-proven testing and interoperability



# Khronos Mission

E.g. OpenGL ES provides 3D on almost **every** smartphone - over 4 Billion units shipped



**Khronos** is an International Industry Consortium of over 100 companies creating royalty-free, **open standard APIs** to enable software to access hardware acceleration for **3D graphics, Virtual and Augmented Reality, Parallel Computing, Neural Networks and Vision Processing**

# Why Mobile and Embedded Need API Standards

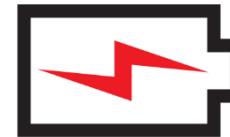
## Without API Standards



Platform  
Fragmentation



Everything  
runs on CPU



## With API Standards



Application  
Portability



Silicon  
Acceleration



**Standard Acceleration APIs provide PERFORMANCE, POWER AND PORTABILITY**

# Khronos Open APIs for Today's Industry



Vision and sensor processing - including neural networks



Interact with sensor, haptic and VR/AR display devices

Download 3D augmentation object and scene data

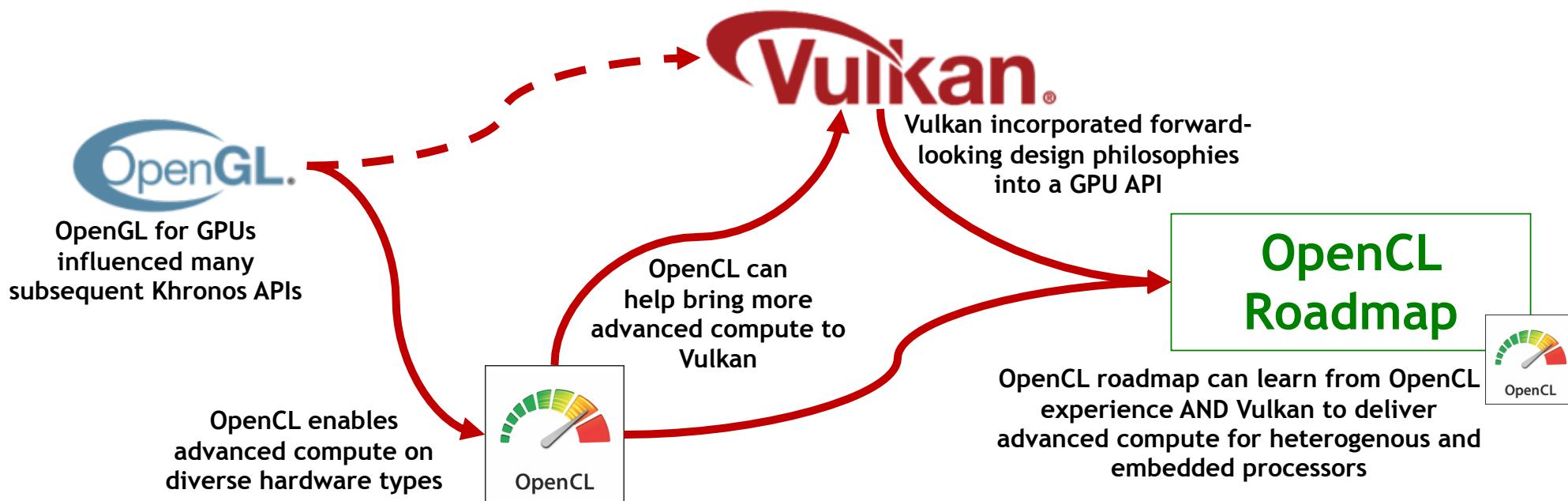


Generate Low Latency 3D Content and Simulations



# Agenda for this Session

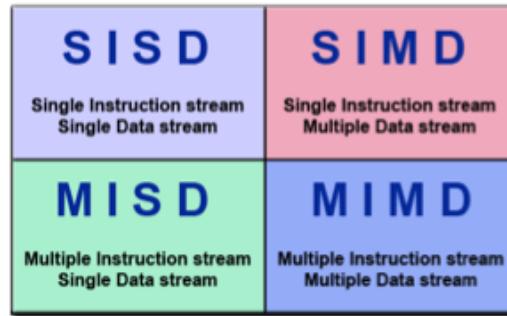
- Introduction to Khronos and the need for open standards
- Update on OpenCL (and SPIR-V) - the standards for heterogenous computing
- Update on Vulkan - a new explicit-style API
- Lessons from Vulkan and OpenCL that can guide the OpenCL roadmap
- Directions for OpenCL roadmap - and increasing relevance to the embedded industry



# Need for Heterogeneous Parallelism



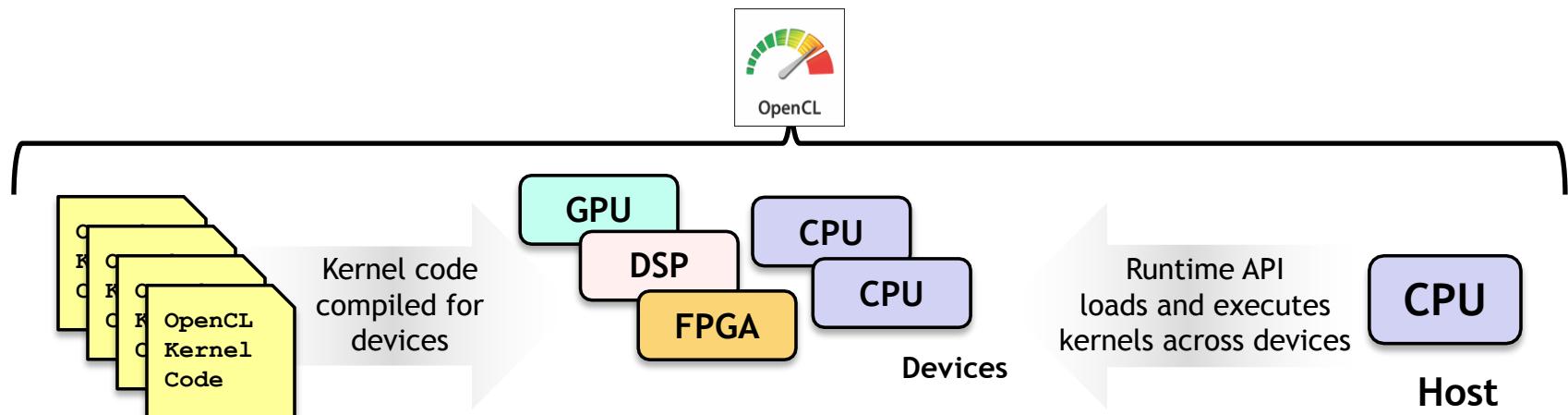
Flynn's Taxonomy (1966)



"The purpose of abstraction is not to be vague, but to create a new semantic level in which one can be absolutely precise" - Edsger Dijkstra

# OpenCL - Low-level Parallel Programming

- Low-level programming of heterogeneous parallel compute resources
  - One code tree can be executed on CPUs, GPUs, DSPs and FPGA ...
- OpenCL C or C++ language to write kernel programs to execute on any compute device
  - Platform Layer API - to query, select and initialize compute devices
  - Runtime API - to build and execute kernels programs on multiple devices
- The programmer gets to control:
  - What programs execute on what device
  - Where data is stored in various speed and size memories in the system
  - When programs are run, and what operations are dependent on earlier operations



# OpenCL 2.2 Released in May 2017



OpenCL 1.2

Becomes industry baseline for heterogeneous parallel computing

OpenCL 2.0

Enables new class of hardware SVM  
Generic Addresses On-device dispatch

OpenCL 2.1  
SPIR-V 1.0

SPIR-V in Core  
Kernel Language  
Flexibility

OpenCL 2.2  
SPIR-V 1.2

OpenCL C++ Kernel Language  
Static subset of C++14  
Templates and Lambdas

SPIR-V 1.2  
OpenCL C++ support

Pipes

Efficient device-scope communication between kernels

## Code Generation Optimizations

- Specialization constants at SPIR-V compilation time
- Constructors and destructors of program scope global objects
- User callbacks can be set at program release time

<https://www.khronos.org/opencl/>

OpenCL 2.2 Reference Card

OpenCL Open Computing Language is a multi-vendor open standard for general-purpose programming of heterogeneous systems that include CPUs, GPUs, and other processors. OpenCL provides a uniform programming environment for fast, efficient, and portable code capable for high-performance compute servers, desktop computer systems, and handheld devices.

Specification documents and online reference are available at [www.khronos.org/opencl](http://www.khronos.org/opencl).

### OpenCL API Reference

#### The OpenCL Platform Layer

The OpenCL platform layer implements platform-specific functionality such as device creation, memory allocation, device configuration information, and to create OpenCL contexts and command queues.

Section and table references are to the OpenCL API 2.2 specification.

`d_int clCreatePlatform(cl_device_id device_id)`  
`d_int clDeletePlatform(cl_device_id device_id)`

**Contents [4]**  
**clCreateContext [4]**  
`cl_int clCreateContext([4]cl_context_properties properties, cl_device_id device_id, void *private_info, size_t private_info_size, cl_int *err_code, void **context)`

**clCreateCommandQueue [4]**  
`cl_int clCreateCommandQueue(cl_device_id device_id, cl_command_queue_properties properties, cl_uint num_queues, cl_int *err_code, cl_command_queue **queue)`

**clCreateCommandQueue [4]**  
`cl_int clCreateCommandQueue(cl_device_id device_id, cl_command_queue_properties properties, cl_uint num_queues, cl_int *err_code, cl_command_queue **queue)`

**clCreateDevice [4]**  
`cl_int clCreateDevice(cl_platform_id platform_id, cl_device_type device_type, cl_device_id *device_id, cl_int *err_code, void **device)`

**clCreateDevice [4]**  
`cl_int clCreateDevice(cl_platform_id platform_id, cl_device_type device_type, cl_device_id *device_id, cl_int *err_code, void **device)`

**clCreateEvent [4]**  
`cl_int clCreateEvent(cl_device_id device_id, cl_int flags, cl_int *err_code, cl_event **event)`

**clCreateEvent [4]**  
`cl_int clCreateEvent(cl_device_id device_id, cl_int flags, cl_int *err_code, cl_event **event)`

**clCreateFence [4]**  
`cl_int clCreateFence(cl_device_id device_id, cl_int flags, cl_int *err_code, cl_fence **fence)`

**clCreateFence [4]**  
`cl_int clCreateFence(cl_device_id device_id, cl_int flags, cl_int *err_code, cl_fence **fence)`

**clCreateHistogram [4]**  
`cl_int clCreateHistogram(cl_device_id device_id, cl_int type, cl_int bins, cl_int *err_code, cl_histogram **histogram)`

**clCreateHistogram [4]**  
`cl_int clCreateHistogram(cl_device_id device_id, cl_int type, cl_int bins, cl_int *err_code, cl_histogram **histogram)`

**clCreateImage [4]**  
`cl_int clCreateImage(cl_device_id device_id, cl_int flags, cl_int type, cl_int width, cl_int height, cl_int depth, cl_int image_base_alignment, cl_int image_max_pitch, cl_int *err_code, cl_image **image)`

**clCreateImage [4]**  
`cl_int clCreateImage(cl_device_id device_id, cl_int flags, cl_int type, cl_int width, cl_int height, cl_int depth, cl_int image_base_alignment, cl_int image_max_pitch, cl_int *err_code, cl_image **image)`

**clCreateMemory [4]**  
`cl_int clCreateMemory(cl_device_id device_id, cl_int flags, cl_int type, cl_int size, cl_int alignment, cl_int *err_code, cl_mem **mem)`

**clCreateMemory [4]**  
`cl_int clCreateMemory(cl_device_id device_id, cl_int flags, cl_int type, cl_int size, cl_int alignment, cl_int *err_code, cl_mem **mem)`

**clCreateProgramWithSource [4]**  
`cl_int clCreateProgramWithSource(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithSource [4]**  
`cl_int clCreateProgramWithSource(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithBinary [4]**  
`cl_int clCreateProgramWithBinary(cl_device_id device_id, cl_int num_binaries, const cl_program_binary *binaries, const cl_program_binary_properties *binary_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithBinary [4]**  
`cl_int clCreateProgramWithBinary(cl_device_id device_id, cl_int num_binaries, const cl_program_binary *binaries, const cl_program_binary_properties *binary_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithBinary [4]**  
`cl_int clCreateProgramWithBinary(cl_device_id device_id, cl_int num_binaries, const cl_program_binary *binaries, const cl_program_binary_properties *binary_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithBinary [4]**  
`cl_int clCreateProgramWithBinary(cl_device_id device_id, cl_int num_binaries, const cl_program_binary *binaries, const cl_program_binary_properties *binary_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithSource [4]**  
`cl_int clCreateProgramWithSource(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithSource [4]**  
`cl_int clCreateProgramWithSource(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithBinary [4]**  
`cl_int clCreateProgramWithBinary(cl_device_id device_id, cl_int num_binaries, const cl_program_binary *binaries, const cl_program_binary_properties *binary_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithBinary [4]**  
`cl_int clCreateProgramWithBinary(cl_device_id device_id, cl_int num_binaries, const cl_program_binary *binaries, const cl_program_binary_properties *binary_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

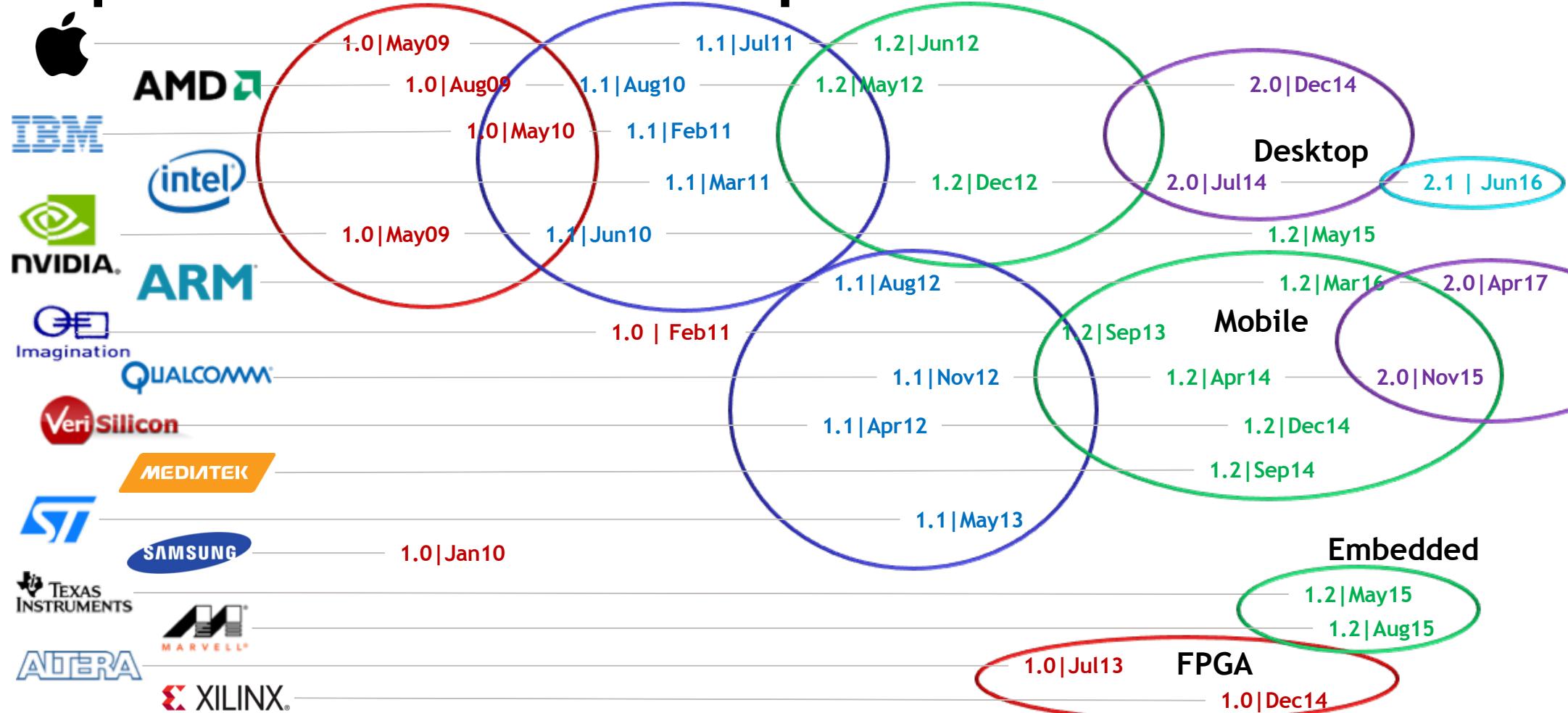
**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

**clCreateProgramWithProperties [4]**  
`cl_int clCreateProgramWithProperties(cl_device_id device_id, cl_int num_sources, const cl_program_source *sources, const cl_program_source_properties *source_properties, cl_int *err_code, cl_program **program)`

# OpenCL Conformant Implementations



Vendor timelines are  
first conformant  
submission for each spec  
generation

Dec08  
OpenCL 1.0  
Specification

Jun10  
OpenCL 1.1  
Specification

Nov11  
OpenCL 1.2  
Specification

Nov13  
OpenCL 2.0  
Specification

Nov15  
OpenCL 2.1  
Specification

# OpenCL Adoption

- 100s of applications using OpenCL acceleration
  - Rendering, visualization, video editing, simulation, image processing
- Over 4,000 GitHub repositories using OpenCL
  - Tools, applications, libraries, languages
  - Up from 2,000 two years ago
- Multiple silicon and open source implementations
  - Increasingly used for embedded vision and neural network inferencing
- Khronos Resource Hub

<https://www.khronos.org/opencl/resources/opencl-applications-using-opencl>

A screenshot of a GitHub search results page for the query "opencl". The header shows the search bar with "opencl", and metrics for 4K repositories, 1M code, 204K commits, 32K issues, 4K wikis, and 123 users. Below the header, it says "4,310 repository results" and "Sort: Most stars". The first two results are listed:

- arrayfire/arrayfire** (C++) - A general purpose GPU library. Last updated 6 days ago.
- LWJGL/lwjgl3** (Kotlin) - A Java library that enables cross-platform access to popular native APIs useful in the

# OpenCL as Language/Library Backend

Caffe

Halide

C++ AMP  
Accelerated Massive Parallelism  
with Microsoft Visual C++™

SYCL™

aparapi

OpenCV

OpenACC®  
DIRECTIVES FOR ACCELERATORS

TensorFlow

C++ based  
Neural  
network  
framework

Language for  
image  
processing and  
computational  
photography

MulticoreWare  
open source  
project on  
Bitbucket

Single  
Source C++  
Programming  
for OpenCL

Java language  
extensions  
for  
parallelism

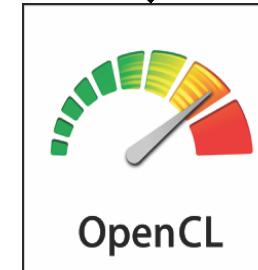
Vision  
processing  
open source  
project

Compiler  
directives for  
Fortran,  
C and C++

Open source  
software library  
for machine  
learning

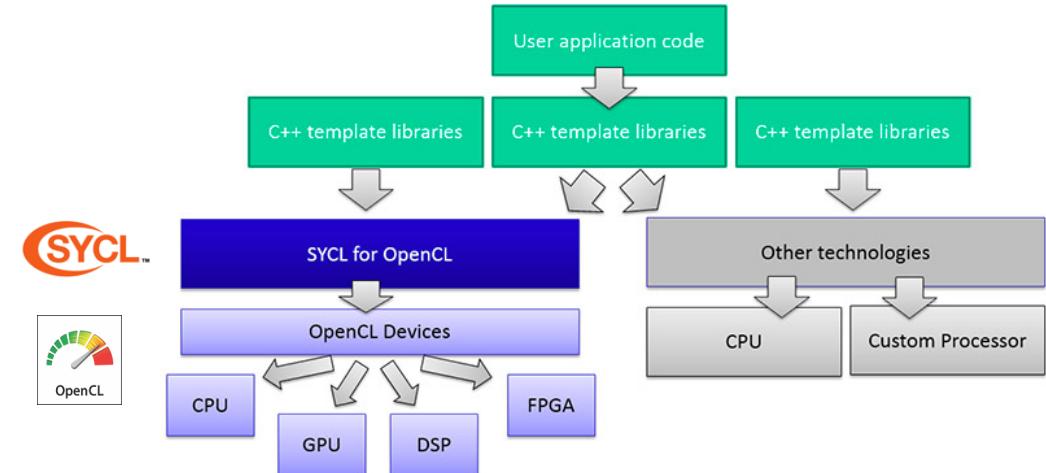
Hundreds of languages, frameworks and projects using OpenCL to access vendor-optimized, heterogeneous compute runtimes

Over 4,000 GitHub repositories using OpenCL: tools, applications, libraries, languages - up from 2,000 two years ago



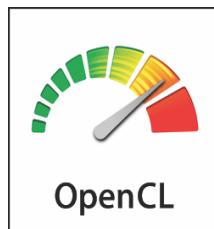
# SYCL Ecosystem

- Single-source heterogeneous programming using STANDARD C++
  - Use C++ templates and lambda functions for host & device code
  - Layered over OpenCL
- Fast and powerful path for bring C++ apps and libraries to OpenCL
  - C++ Kernel Fusion - better performance on complex software than hand-coding
  - Halide, Eigen, Boost.Compute, SYCLBLAS, SYCL Eigen, SYCL TensorFlow, SYCL GTX
  - triSYCL, ComputeCpp, VisionCpp, ComputeCpp SDK ...
- More information at <http://sycl.tech>



# The Choice of SYCL 2.2 or OpenCL C++

**C++ Kernel Language**  
Low Level Control  
'GPGPU'-style separation  
of device-side kernel  
source code and host code



**Developer Choice**  
The development of the two  
specifications are aligned so code  
can be easily shared between the  
two approaches

**Single-source C++**  
Programmer Familiarity  
Approach also taken by  
C++ AMP, OpenMP and the  
C++ 17 Parallel STL



# SPIR-V Transforms the Language Ecosystem

- First multi-API, intermediate language for parallel compute and graphics
  - Natively represents structures in shader and kernel languages
  - <https://www.khronos.org/registry/spir-v/papers/WhitePaper.pdf>
- Compiler IR for OpenCL, Vulkan and OpenGL
  - Easy to parse - just a stream of words
  - Easy to transform - designed to be easy to convert to and from LLVM IR
  - Easy to manipulate and optimize - single-Static Assignment form

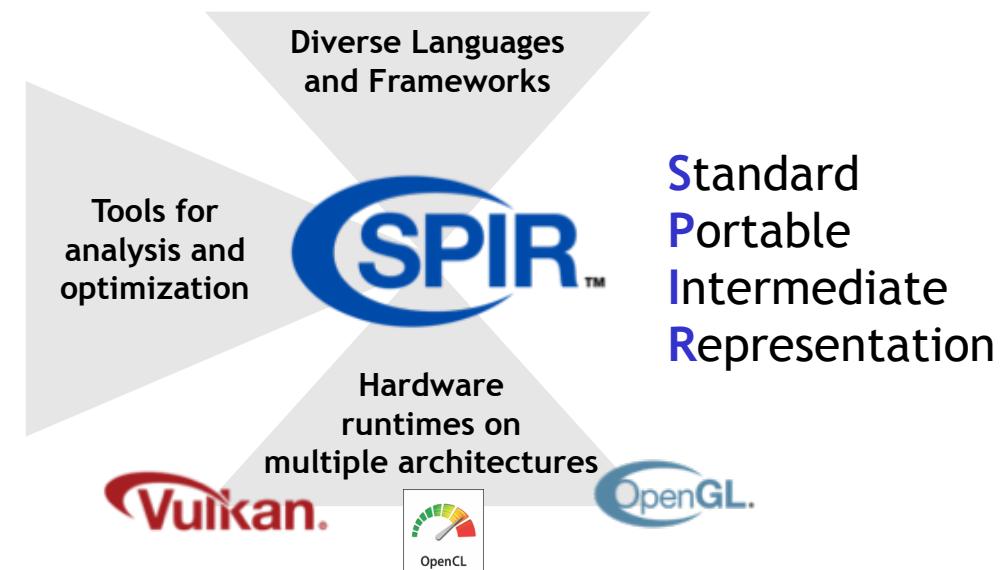
## Multiple Developer Advantages

Use same front-end compiler for all platforms

Ship SPIR-V - not shader source code

Simpler and more reliable drivers

Reduces runtime kernel compilation time

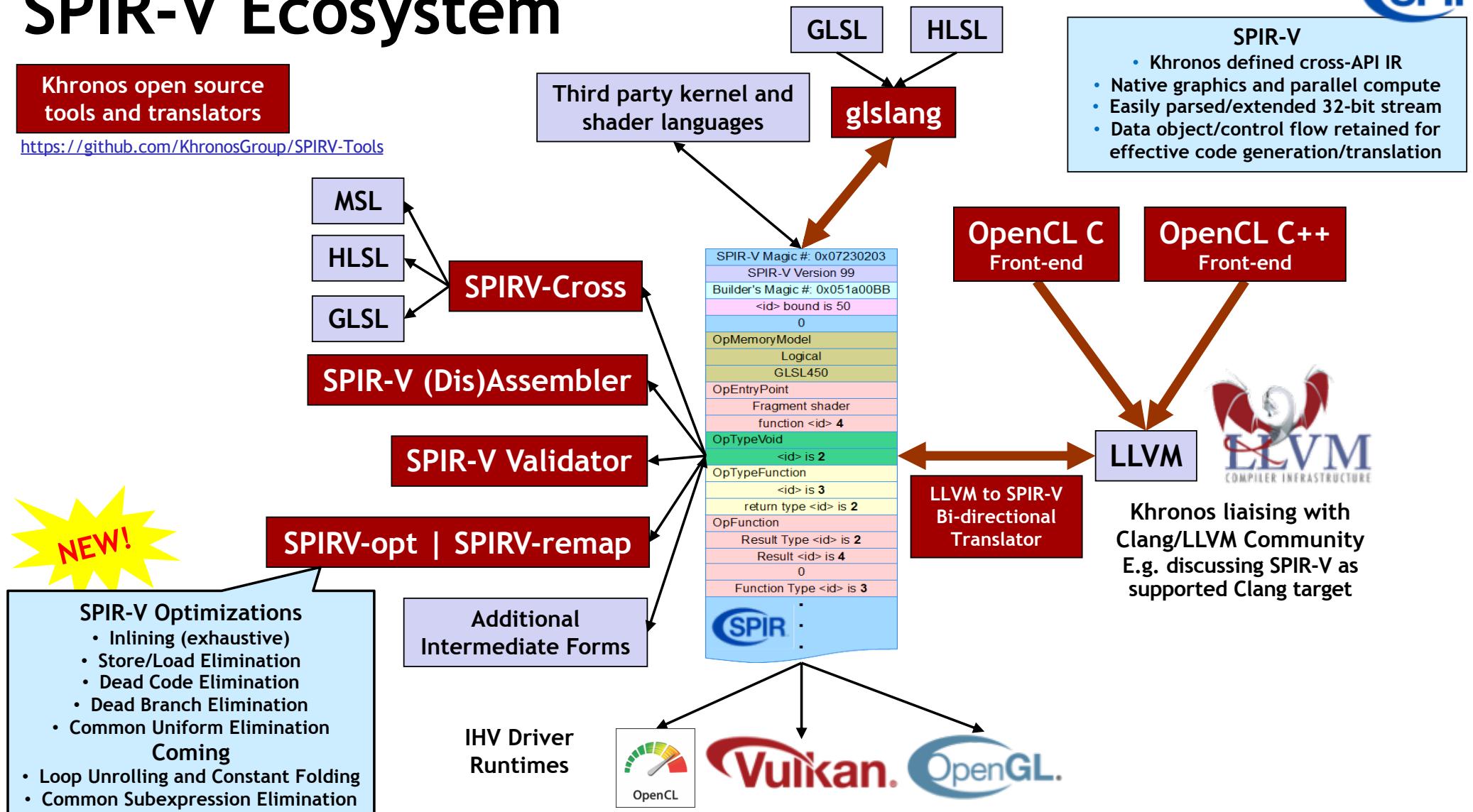


# SPIR-V Ecosystem



Khronos open source tools and translators

<https://github.com/KhronosGroup/SPIRV-Tools>



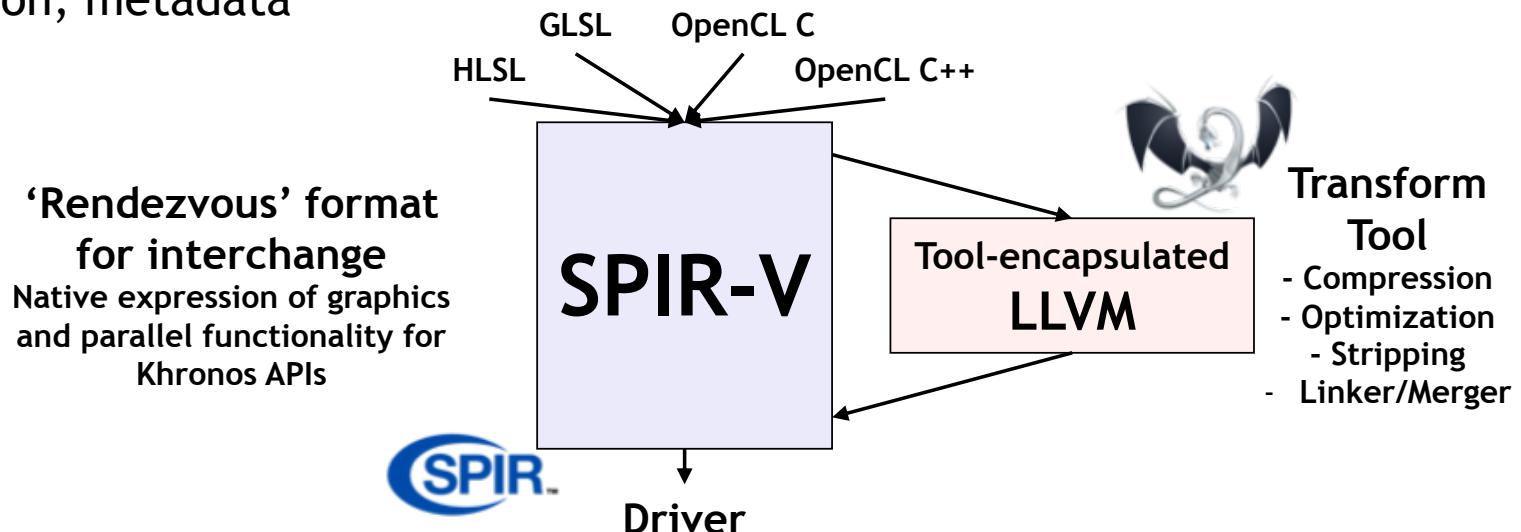
# Evolution of SPIR Family

- SPIR-V is first fully specified Khronos-defined SPIR standard
  - Does not use LLVM to isolate from LLVM roadmap changes
  - Includes full flow control, graphics and parallel constructs beyond LLVM
  - Khronos open sourcing SPIR-V <-> LLVM conversion tools

	SPIR 1.2	SPIR 2.0	SPIR-V 1.X
LLVM Interaction	Uses LLVM 3.2	Uses LLVM 3.4	100% Khronos defined Round-trip lossless conversion
Compute Constructs	Metadata/Intrinsics	Metadata/Intrinsics	Native
Graphics Constructs	No	No	Native
Supported Language Feature Set	OpenCL C 1.2	OpenCL C 1.2 OpenCL C 2.0	OpenCL C 1.2 / 2.X OpenCL C++ GLSL
OpenCL Ingestion	OpenCL 1.2 Extension	OpenCL 2.0 Extension	OpenCL 2.1/2.2 Core
Graphics API Ingestion	-	-	Vulkan and OpenGL 4.6 Core

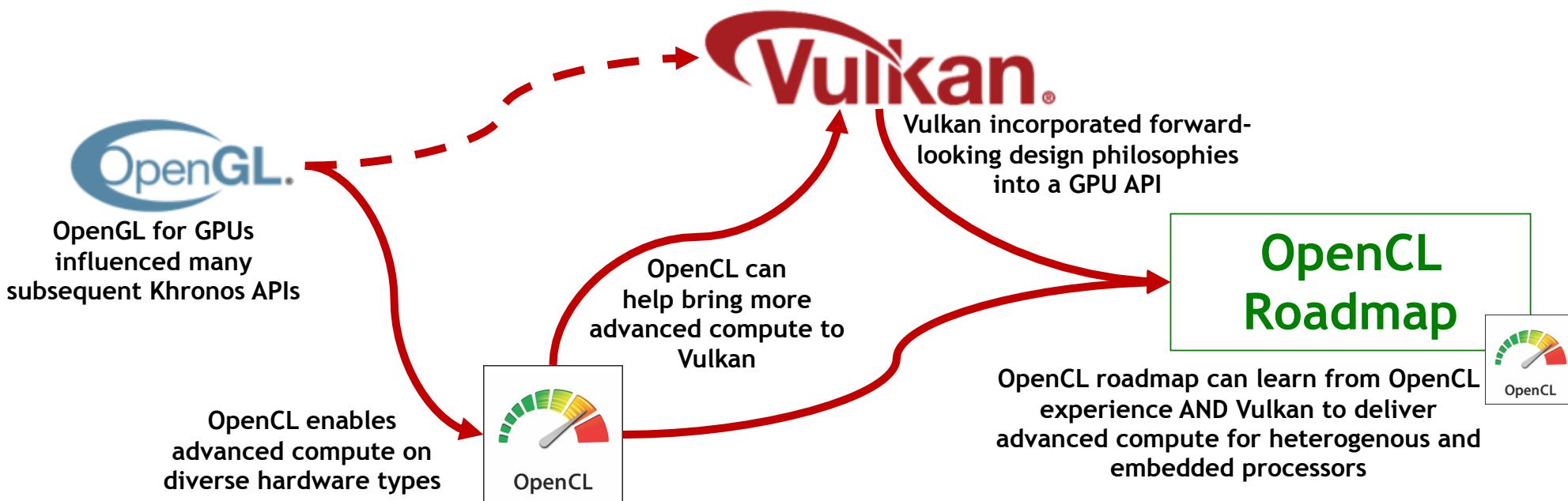
# Support for Both SPIR-V and LLVM

- LLVM is an SDK, not a formally defined standard
  - Khronos moved away from trying to use LLVM IR as a standard
  - Issues with versioning, metadata, etc.
- But LLVM is a treasure chest of useful transforms
  - SPIR-V tools can encapsulation and use LLVM to do useful SPIR-V transforms
- SPIR-V tools can all use different rules - and there will be lots of these
  - May be lossy and only support SPIR-V subset
  - Internal form is not standardized
  - May hide LLVM version, metadata



# Agenda for this Session

- Introduction to Khronos and the need for open standards
- Update on OpenCL (and SPIR-V) - the standards for heterogenous computing
- **Update on Vulkan - a new explicit-style API**
- Lessons from Vulkan and OpenCL that can guide the OpenCL roadmap
- Directions for OpenCL roadmap - and increasing relevance to the embedded industry



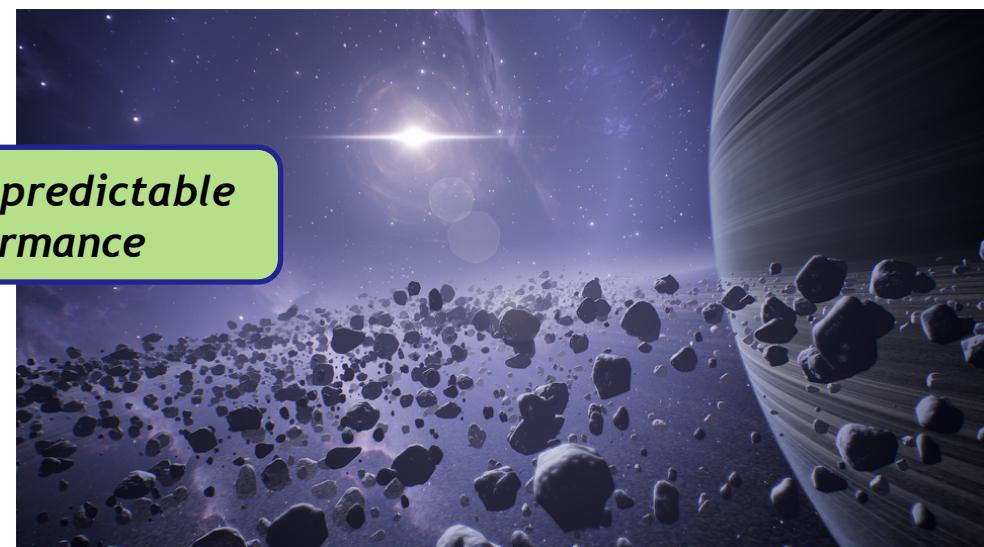
# The Key Principle of Vulkan: *Explicit Control*

- Application tells the driver what it is going to do
  - *In enough detail* that driver doesn't have to guess
  - *When* the driver needs to know it
- In return, driver promises to do
  - *What* the application asks for
  - *When* it asks for it
  - *Very quickly*
- *No driver magic - no surprises*

You have to supply a  
LOT of information!

And, you have to  
supply it early

*Efficient, predictable  
performance*

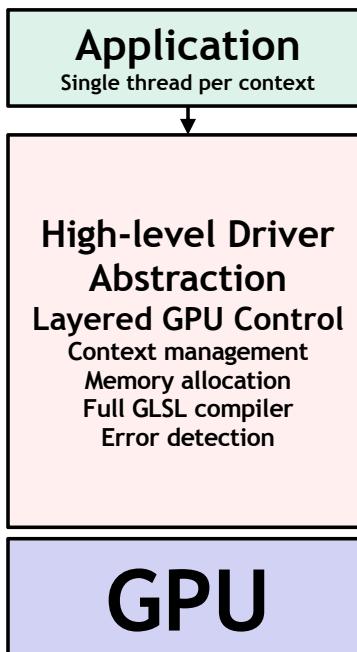


# Vulkan Explicit GPU Control

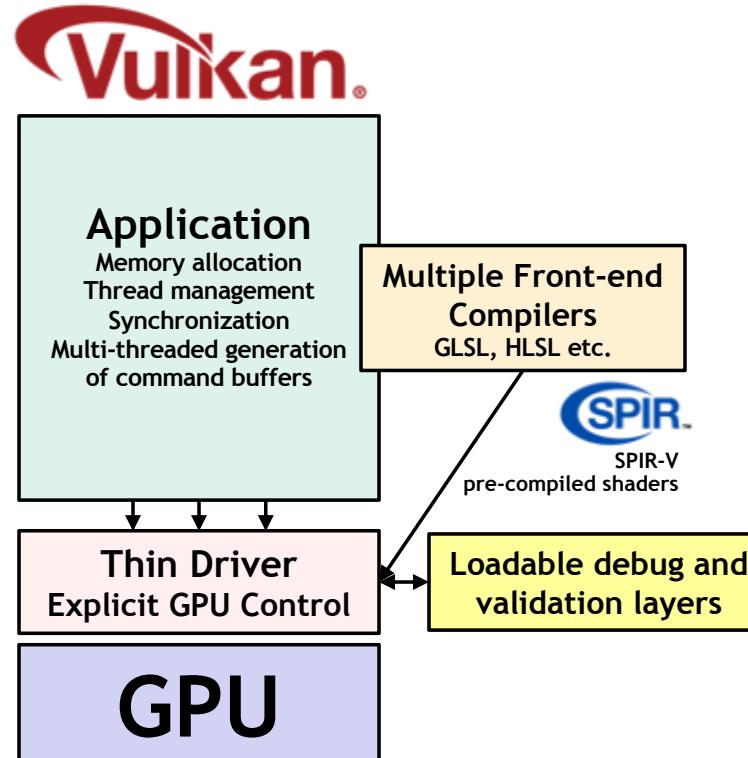
Vulkan = high performance and low latency 3D and GPU Compute. Ideal for VR/AR applications



Complex drivers cause overhead and inconsistent behavior across vendors  
Always active error handling  
Full GLSL preprocessor and compiler in driver  
OpenGL vs. OpenGL ES



Vulkan 1.0 provides access to OpenGL ES 3.1 / OpenGL 4.X-class GPU functionality but with increased performance and flexibility



Resource management offloaded to app: low-overhead, low-latency driver  
Consistent behavior: no ‘fighting with driver heuristics’  
Validation and debug layers loaded only when needed  
SPIR-V intermediate language: shading language flexibility  
Multi-threaded command creation. Multiple graphics, command and DMA queues  
Unified API across all platforms with feature set flexibility

# The Need for Cross-Platform 3D Standards

## Vulkan Provides

Clean, modern architecture | Low overhead, explicit GPU access  
Portable across desktop and mobile | Multi-thread / multi-core friendly  
Primary 3D API on Android Platform  
*Efficient, low-latency, predictable performance*  
*Access to multiple platforms*



macOS

# Pervasive Vulkan



All Major GPU Companies shipping Vulkan Drivers for Desktop and Mobile Platforms



<http://vulkan.gpuinfo.org/>

## Desktop, Mobile, Embedded and Console Platforms Supporting Vulkan

Including phones and tablets from Google, Huawei, Samsung, Sony, Xiaomi - both premium and mid-range devices



Desktop



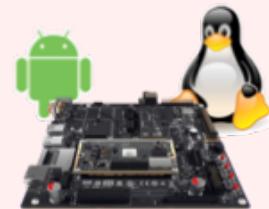
Android 7.0+



Nintendo Switch



Android TV



Embedded Linux



SteamVR



GearVR



Oculus Rift



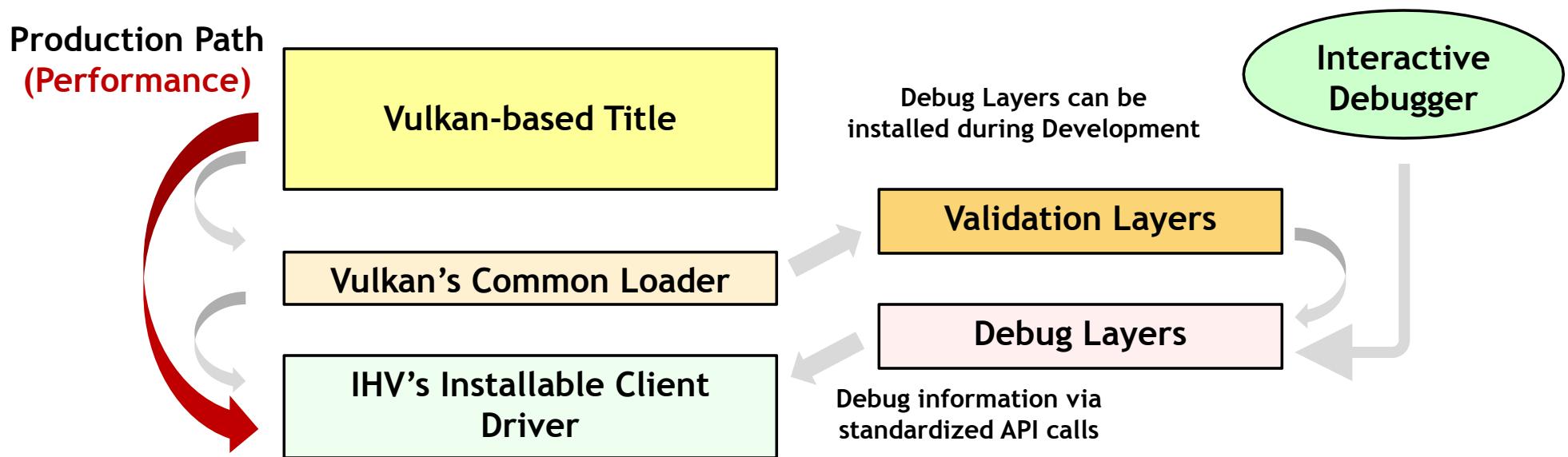
Google Daydream

## Game Engines



# Vulkan Tools Architecture

- Layered design for cross-vendor tools innovation and flexibility
  - IHVs plug into a common, extensible architecture for code validation, debugging and profiling during development without impacting production performance
- Khronos Open Source Loader enables use of tools layers during debug
  - Finds and loads drivers, dispatches API calls to correct driver and layers



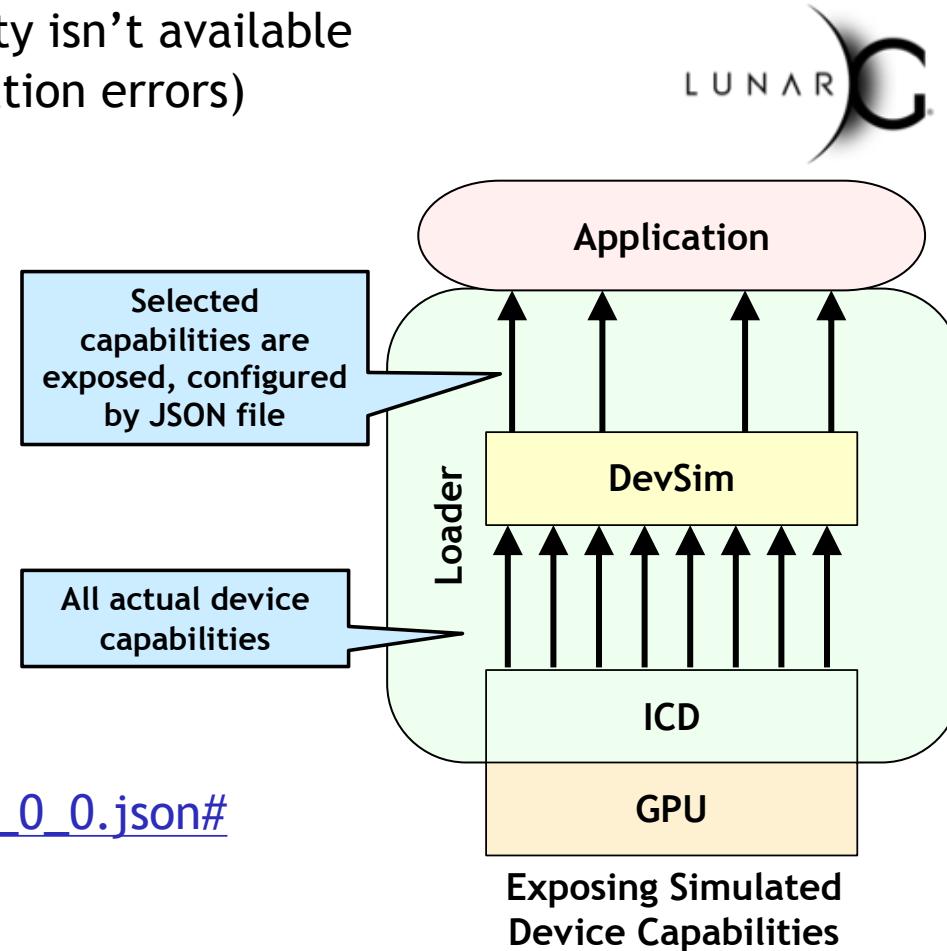
# Vulkan Feature Sets

- Vulkan supports hardware with a wide range of hardware capabilities
  - Mobile OpenGL ES 3.1 up to desktop OpenGL 4.5 and beyond
- One unified API framework for desktop, mobile, console, and embedded
  - No "Vulkan ES" or "Vulkan Desktop"
- Vulkan precisely defines a set of "fine-grained features"
  - Features are specifically enabled at device creation time (similar to extensions)
- Platform owners define a Feature Set for their platform
  - Vulkan provides the mechanism but does not mandate policy
  - Khronos will define Feature Sets for platforms where owner is not engaged
- Khronos will define feature sets for Windows and Linux
  - After initial developer feedback



# VK\_LAYER\_LUNARG\_device\_simulation

- Simulate capabilities of mobile and embedded devices
  - Test application without requiring physical device
  - Exercise fall-back code paths, when a capability isn't available
  - Find unintentional assumptions (triggers validation errors)
- Modifies results from Vulkan queries
  - Device configuration defined by JSON file
  - Integrated with Sascha Willems database
- Simulation, NOT Emulation
  - Doesn't add more capabilities not already present in actual device
- Source available now
  - <https://github.com/LunarG/VulkanTools>
  - Please submit issues
- Verify configuration files are correct
  - [https://schema.khronos.org/vulkan/devsim\\_1\\_0\\_0.json#](https://schema.khronos.org/vulkan/devsim_1_0_0.json#)



# Vulkan Genesis

Khronos members from all segments of the graphics industry agree the need for new generation cross-platform GPU API

Unprecedented level of participation from game engine developers

Significant proposals, IP contributions and engineering effort from many working group members

18 months  
A high-energy working group effort

Khronos' first API 'hard launch' February 2016

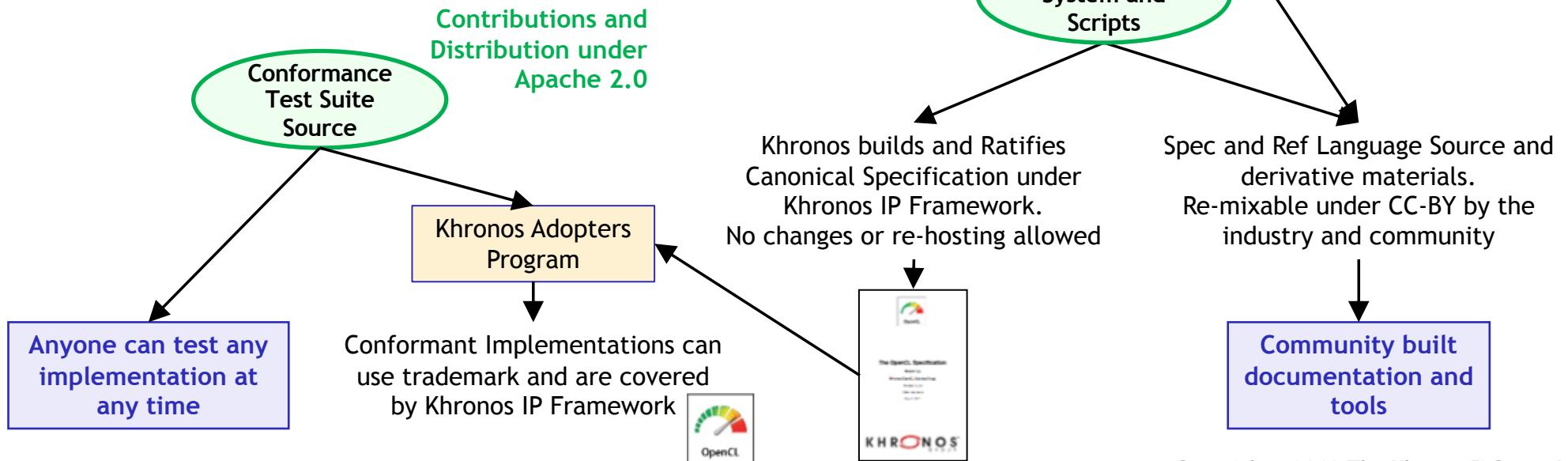
Specifications, Conformance Tests, SDKs, Reference Materials, Compiler front-ends, Samples - all open source...  
Multiple Conformant Drivers on multiple OS



Vulkan Working Group Participants

# New Open Source Engagement Model

- Khronos is open sourcing specification sources, conformance tests, tools
  - Merge requests welcome from the community (subject to review by OpenCL working group)
- Deeper Community Enablement
  - Mix your own documentation!
  - Contribute and fix conformance tests
  - Fix the specification, headers, ICD etc.
  - Contribute new features (carefully)



# Vulkan Evolution



Feb16  
Vulkan 1.0

Explicit Access to GPU Acceleration

## Strengthened Ecosystem and SDK

Enhanced developer and debugging tools

Regression testing for SDK stability

Enhanced Conformance Testing (API now has 198K test cases  
- up from 107K last year)

Compiler robustness - including HLSL support

## Vulkan Extensions

Maintenance updates plus additional functionality

Explicit Building Blocks for VR

Explicit Building Blocks for Homogeneous Multi-GPU

Enhanced Windows System Integration

Increased Shader Language Flexibility

Enhanced Cross-Process and Cross-API Sharing

## Widened Platform Support

Through Vulkan Portability Initiative  
Including Vulkan on macOS and iOS

## Vulkan Roadmap

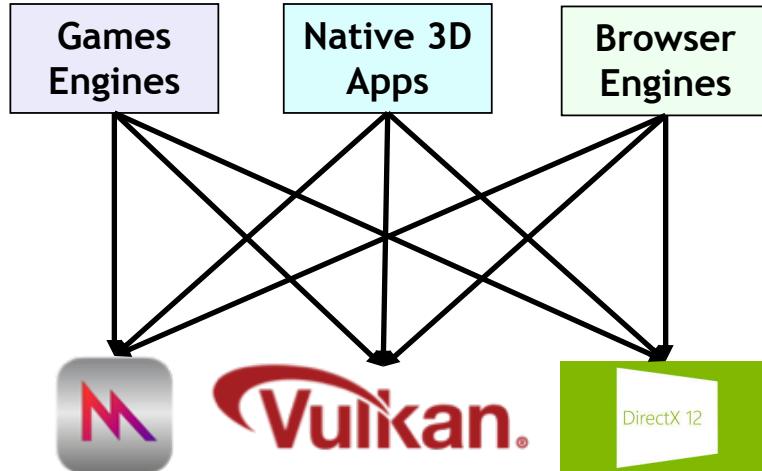
### Regular Vulkan Core Releases

Integrates proven KHR extensions

### Roadmap Discussions

Pushes forward the envelope of GPU Acceleration:  
Enhanced Compute and Language Flexibility  
Optimized Vision and Inferencing Acceleration  
Ray Tracing

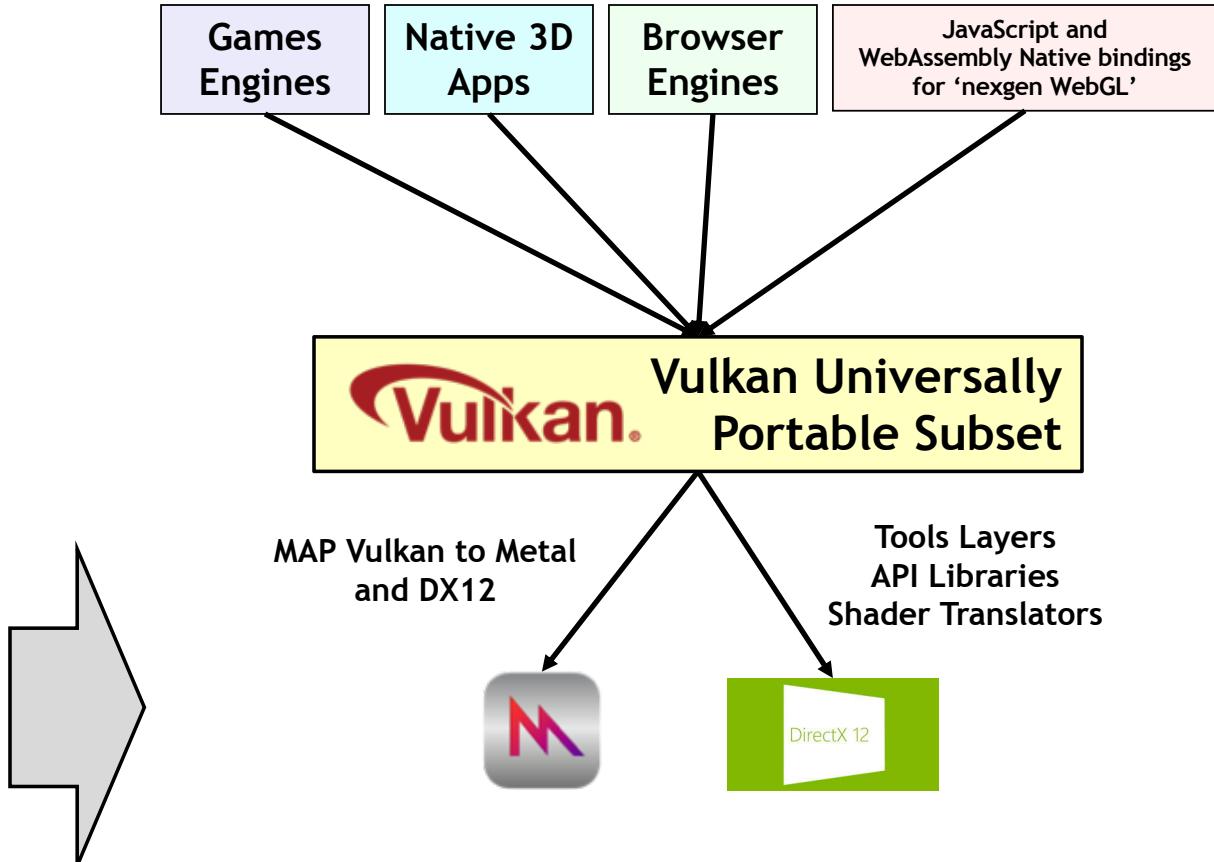
# Market Demand for Universal 3D Portability



Community Outreach at GDC 2017  
Create a hybrid Portability API?

Feedback - AVOID CREATING A FOURTH API!!!

Would need new specification, CTS, Documentation.  
Additional developer learning curve.  
A whole new specification to name, brand, promote.  
Would INCREASE industry fragmentation



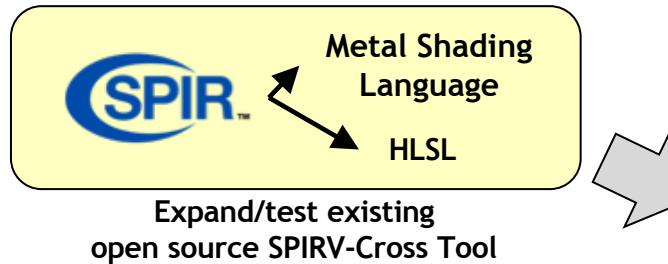
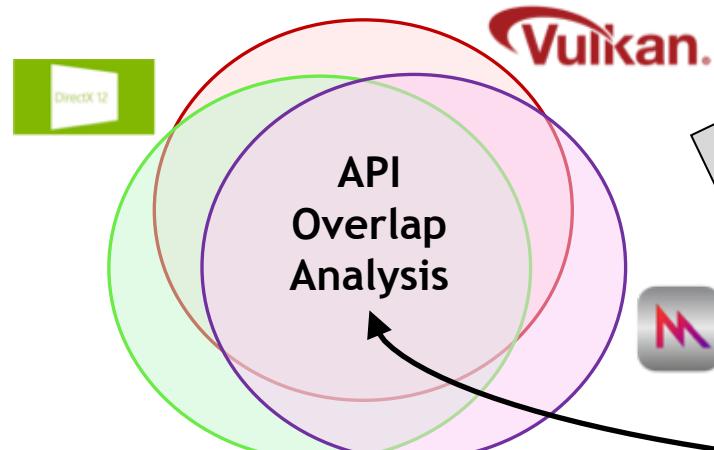
# Vulkan Universal Portability



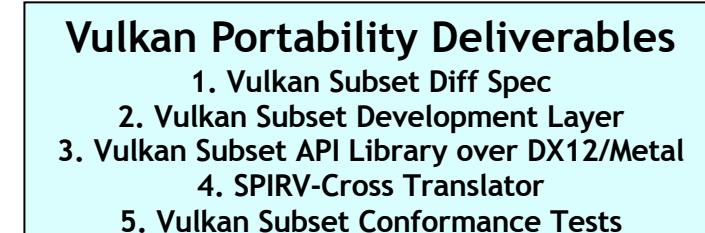
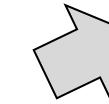
Open source project with similar goals  
<https://github.com/gfx-rs/gfx>



Vulkan on iOS and macOS  
<https://moltengl.com/moltenvk/>

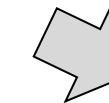


Identify Vulkan features not directly mappable to DX12 and Metal



Layers, APIs, Translators and Tests all to be developed and released in open source

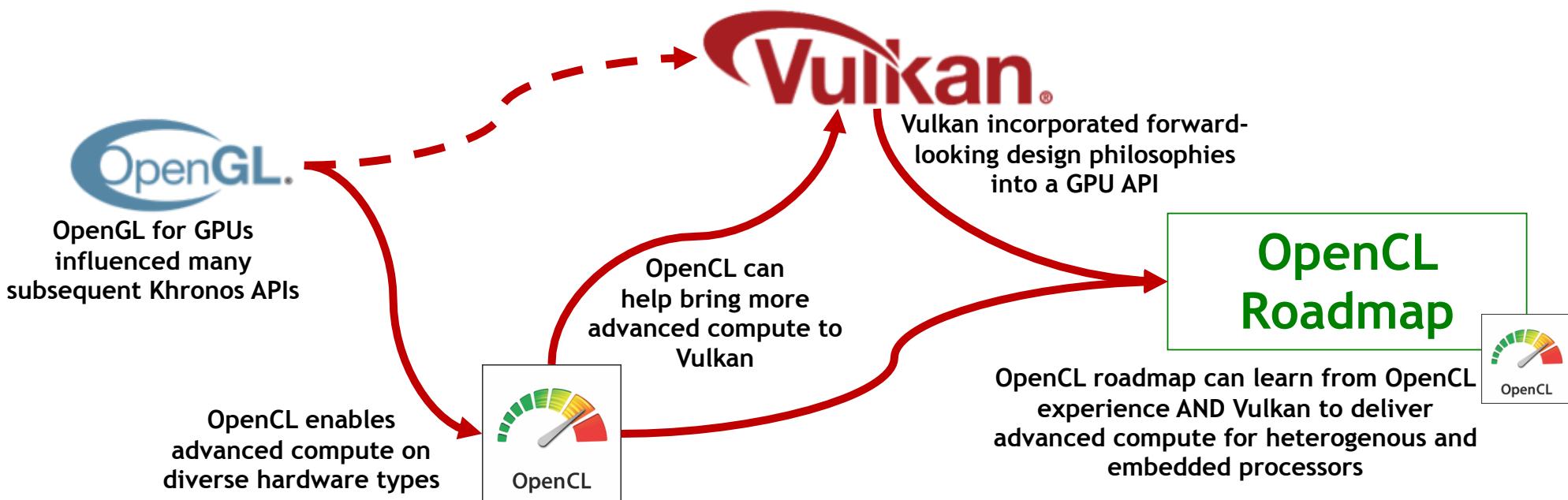
Possible proposals for Vulkan extensions for enhanced portability (and possibly Web robustness) sent to Vulkan WG



New Vulkan functionality may affect the overlap analysis

# Agenda for this Session

- Introduction to Khronos and the need for open standards
- Update on OpenCL (and SPIR-V) - the standards for heterogenous computing
- Update on Vulkan - a new explicit-style API
- **Lessons from Vulkan and OpenCL that can guide the OpenCL roadmap**
- Directions for OpenCL roadmap - and increasing relevance to the embedded industry



# Key Vulkan Lessons

- Engine developer insights were essential during design
  - Engine prototyping during design was essential during design
- Open sourcing tests, tools, specs drives deeper community engagement
  - Support strong middleware ecosystem
- Explicit API
  - Provide straightforward hardware access
  - Aim for as much functionality portability as possible



Rely on middleware engines and libraries for

1. Ease of Use
2. Performance Portability

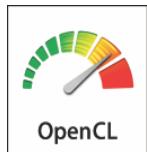


# OpenCL - 1000s Man Years Effort



## Single Source C++ Programming

Full support for features in C++14-based Kernel Language



## API and Language Specs

Brings C++14-based Kernel Language into core specification



## Portable Kernel Intermediate Language

Support for C++14-based kernel language e.g.  
constructors/destructors

3-component vectors  
Additional image formats  
Multiple hosts and devices  
Buffer region operations  
Enhanced event-driven execution  
Additional OpenCL C built-ins  
Improved OpenGL data/event interop

Device partitioning  
Separate compilation and linking  
Enhanced image support  
Built-in kernels / custom devices  
Enhanced DX and OpenGL Interop

Shared Virtual Memory  
On-device dispatch  
Generic Address Space  
Enhanced Image Support  
C11 Atomics  
Pipes  
Android ICD

SPIR-V in Core  
Subgroups into core  
Subgroup query operations  
clCloneKernel  
Low-latency device timer queries

## OpenCL C++ Kernel Language

Static subset of C++14

Templates and Lambdas

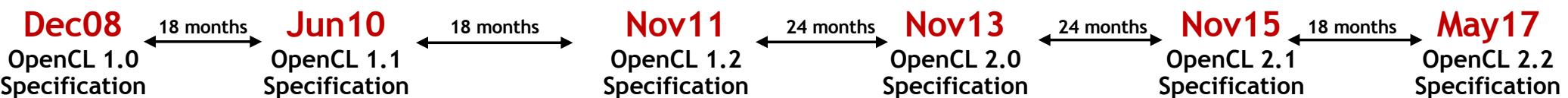
## SPIR-V 1.2 with C++ support

SYCL 2.2 single source C++  
Pipes

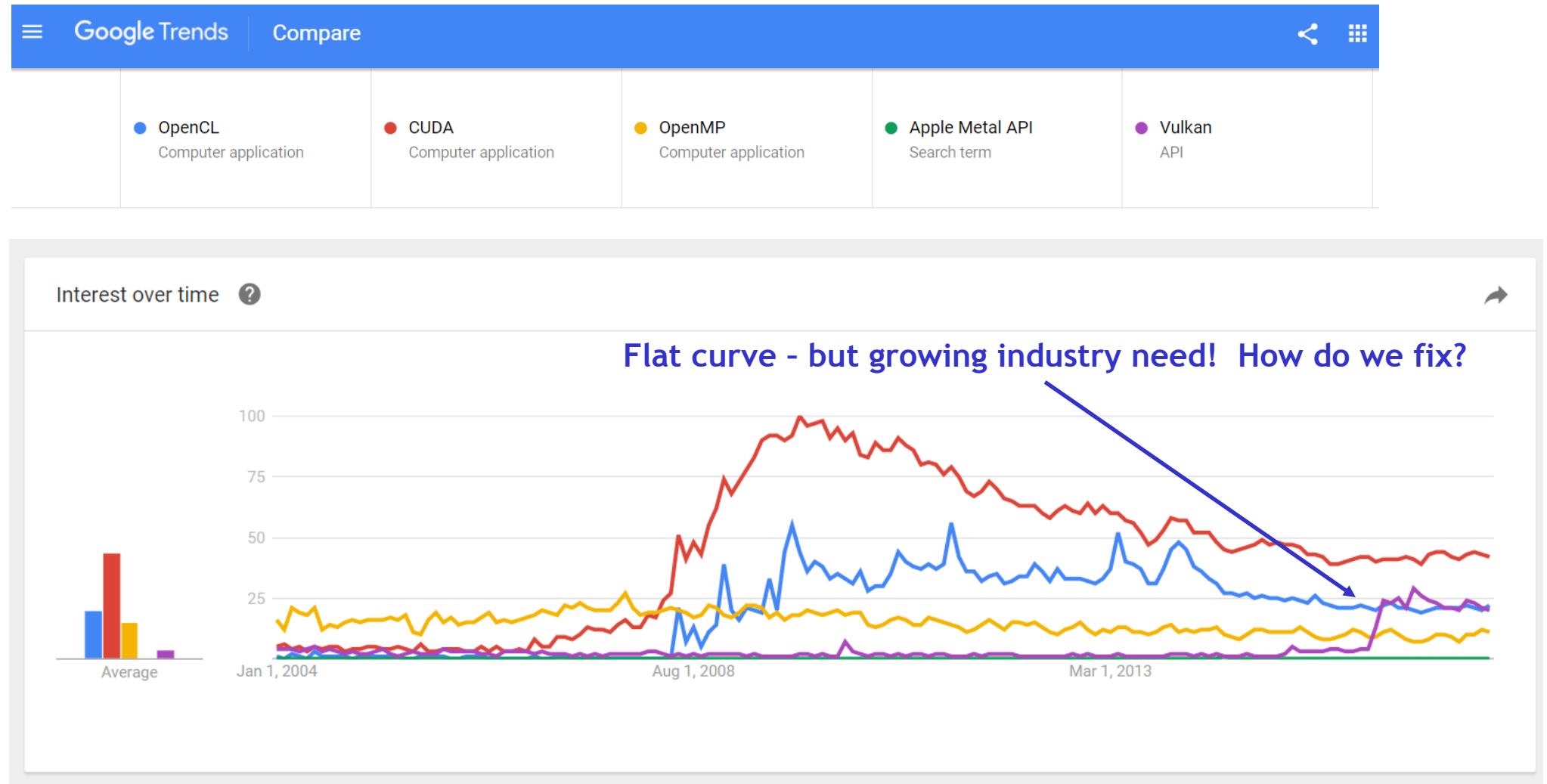
Efficient device-scope  
communication between kernels

## Multiple Code Generation Optimizations

LEARN



# Google Trends



# Key Lessons

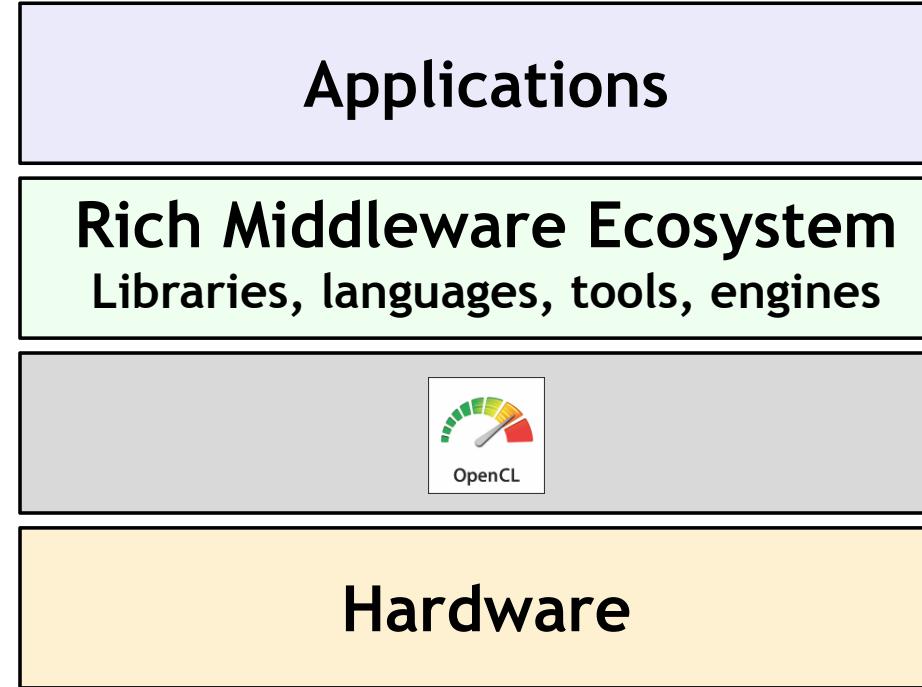
Lessons	How We Learned Them	How We Do Better!
Language flexibility is good! Enable language innovation!	OpenCL WG spent too long designing OpenCL C and C++	Ingest SPIR-V! and open source front-ends
Optimized tools and libraries are key	Lack of highly-optimized libraries for key use cases compared to proprietary solutions	Engage with library ecosystem and ensure the standard enables effective library deployment
Platform adoption is key influencer on overall ecosystem momentum	Apple are focused on Metal RenderScript confusion on Android NVIDIA not pushing to full 2.0	Enable implementations to ship that are relevant to a platform= and its customers
Software developer insights and prototyping are essential during standards design	OpenCL roadmap has been to feature oriented - not fully utilized by developers	Encourage ISVs to join Khronos as members or on OpenCL Advisory Panels

# Embrace the Layered Ecosystem

OpenCL mixes providing low-level hardware access with ‘ease-of-use’

Didn’t make it clear that low-level performance portability is impossible

Did not focus on rapidly porting efficient libraries



Middleware just needs direct access to hardware. Driver should ‘get out of the way’

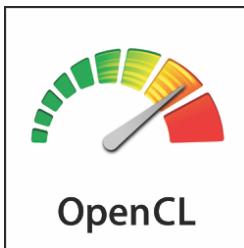
Middleware can provide ease of use

Middleware has the system/domain context to try to provide performance portability

Run-time abstraction to compute hardware IS needed:

- Software vendors can’t afford to port to every type/generation hardware
- Hardware vendors want to keep innovating under an abstraction

# Market Segments Need Deployment Flexibility



OpenCL has been over-monolithic

E.g. DSP inferencing should not be forced to ship IEEE FP32

Solution: feature sets - enabling toggling capabilities within a coherent framework without losing conformance

## Desktop (physical and cloud)

Use cases: Video Image Processing, Gaming Compute, Rendering, Neural Network Training and Inferencing

Roadmap: Vulkan interop, dialable precision, pre-emption, collective programming and improved execution model, dynamic parallelism, pre-emption

## Mobile

Use case: Photo and Vision Processing, Neural Network Inferencing

Roadmap: SVM, dialable precision for inference engine and pixel processing efficiency, pre-emption and QoS scheduling for power efficiency

## HPC

Use case: Numerical Simulation, Neural Network Training, Virtualization

Roadmap: enhanced streaming processing, enhanced library support

## FPGAs

Use cases: Network and Stream Processing

Roadmap: enhanced execution model, self-synchronized and self-scheduled graphs, fine-grained synchronization between kernels, DSL in C++

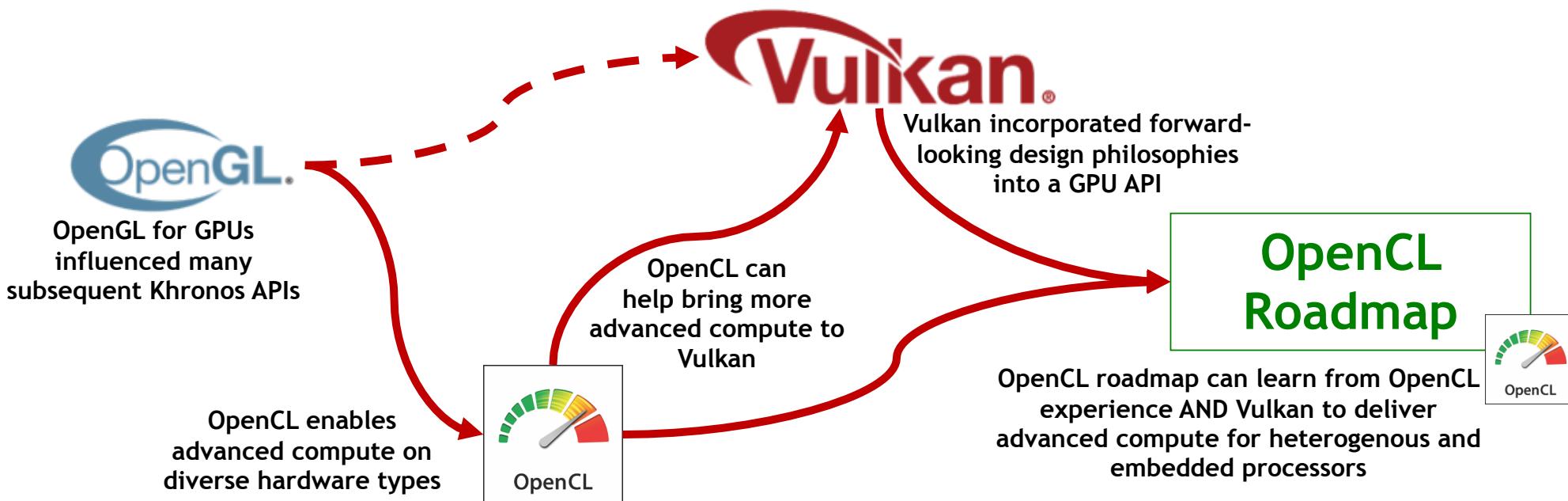
## Embedded

Use cases: Vision, Signal and Pixel Processing, Neural Network and Inferencing

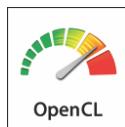
Roadmap: arbitrary precision for power efficiency, hard real-time scheduling, asynch DMA

# Agenda for this Session

- Introduction to Khronos and the need for open standards
- Update on OpenCL (and SPIR-V) - the standards for heterogenous computing
- Update on Vulkan - a new explicit-style API
- Lessons from Vulkan and OpenCL that can guide the OpenCL roadmap
- Directions for OpenCL roadmap - and increasing relevance to the embedded industry

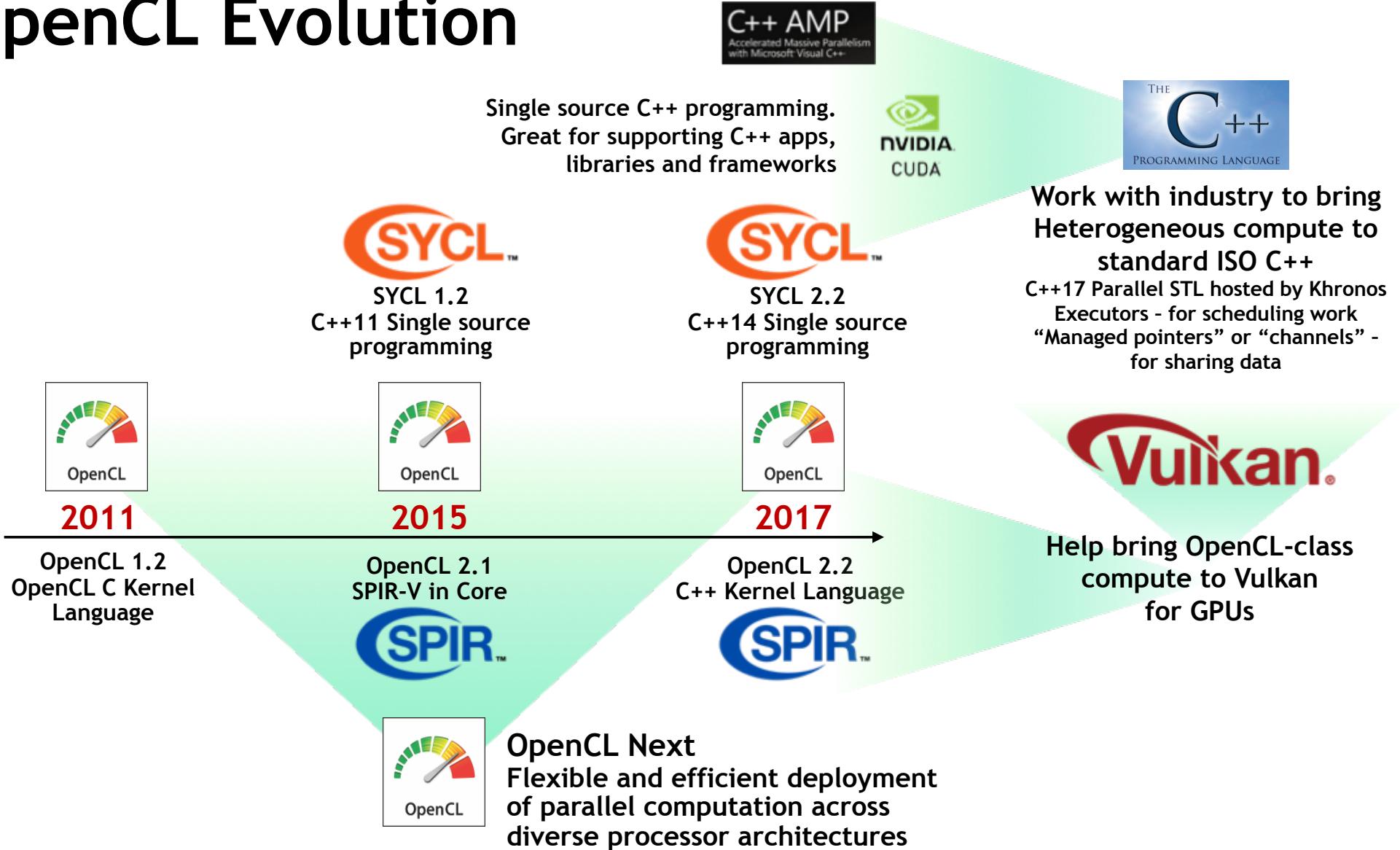


# Requirements for OpenCL Roadmap



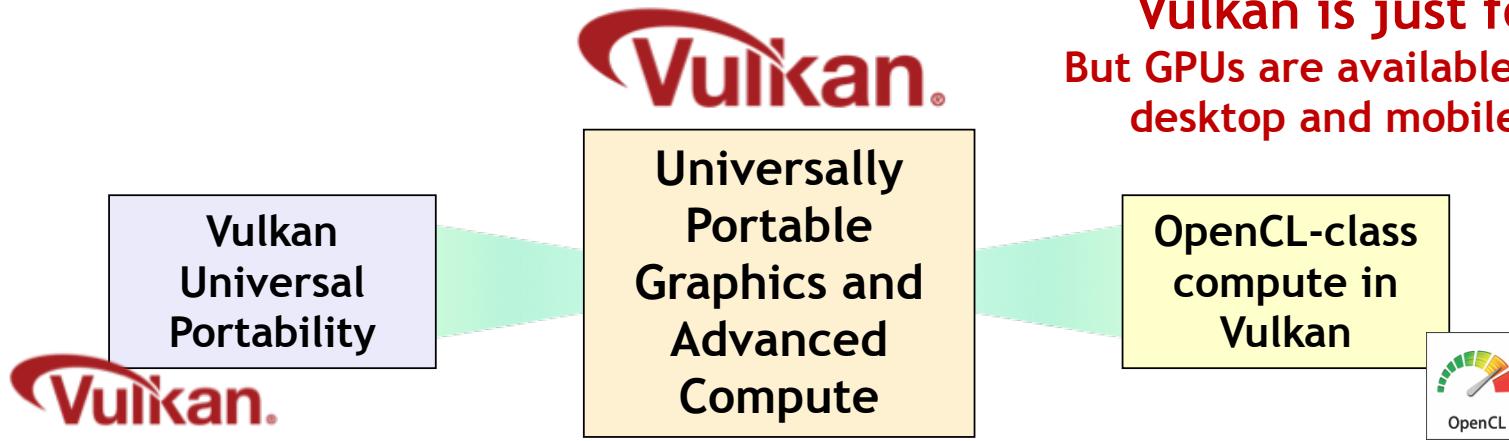
Low-level explicit API as Foundation of multi-layer compute ecosystem	✓	✓
Flexible features selection for deployment flexibility	✓	✓
SPIR-V Ingestion for Language flexibility	✓	✓
Widely Adopted No market barriers to deployment	✓	✓
Installable tools architecture for development and deployment flexibility	✓	✓
Low-latency, multi-threaded dispatch for fine-grained, high-performance	✓	✓
Enables at least OpenCL 2.X-class compute capabilities	✓	✗
Support for diverse processor types	✓	✗

# OpenCL Evolution



# Bringing Advanced Compute to Vulkan

- Already developed Vulkan/SPIR-V functionality to support OpenCL C compute operations
  - 16-bit storage and use of fp16 and int8 in shaders
  - Variable Pointers and Subgroups
- In future will target additional advanced use cases
  - Improve Compute for Graphics (advocated in a number of HPG/SIGGRAPH 2016/17 talks, including Andrew Lauritzen's talk @ Open Problems in Real-Time Rendering, SIGGRAPH'17)
  - Improve support for processing images acquired from camera sensors
  - Enhanced support for performant and power-efficient machine learning

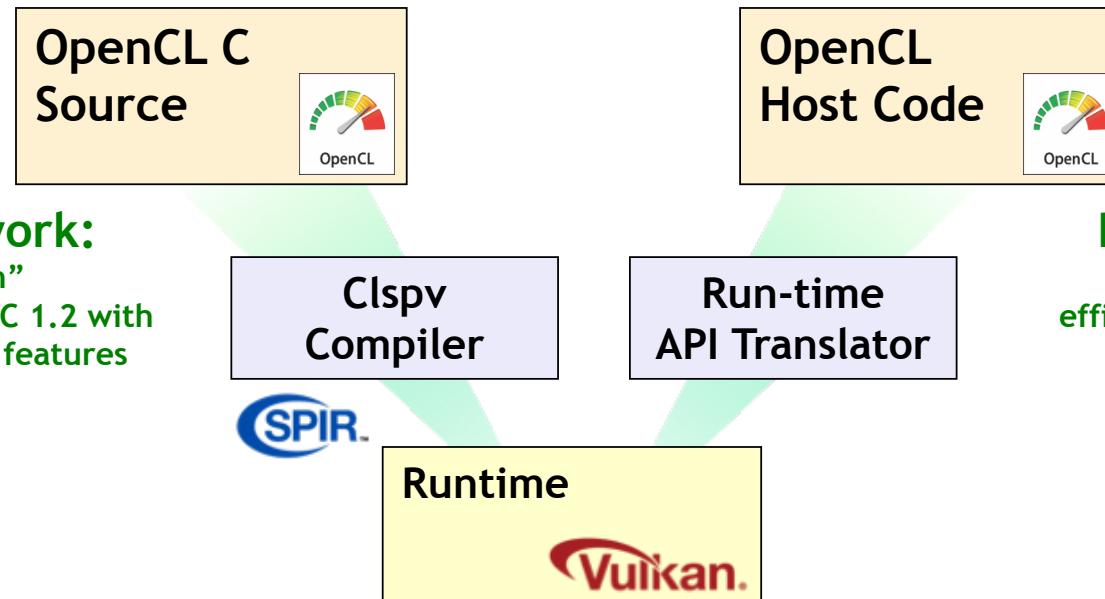


**Vulkan is just for GPUs!**  
But GPUs are available across cloud,  
desktop and mobile platforms

# Clspv OpenCL C to Vulkan Compiler



- Experimental collaboration between Google, Codeplay, and Adobe
  - Successfully tested on over 200K lines of Adobe OpenCL C production code
  - Released in open source <https://github.com/google/clspv>
  - Tracks top-of-tree LLVM and clang, not a fork
- Compiles OpenCL C's programming model to Vulkan's SPIR-V execution environment
  - Proof-of-concept that OpenCL compute can be brought seamlessly to Vulkan



## Possible future work:

“OpenCL C Vulkan Edition”  
Proper subset of OpenCL C 1.2 with  
support for select Vulkan features

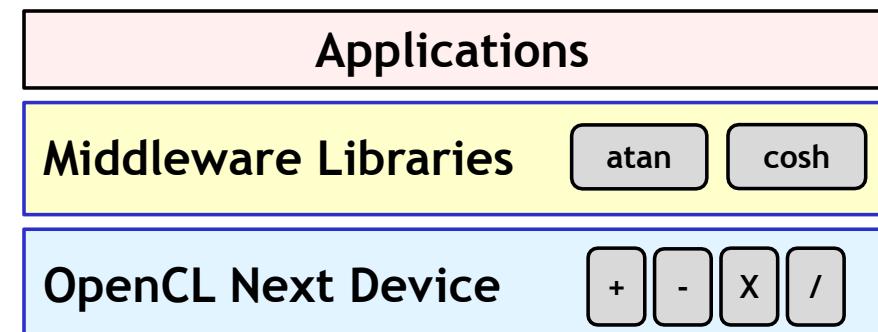
## Possible future work:

Subset of OpenCL API that  
efficiently maps to Vulkan through  
a run-time shim

All tools, compilers  
and conformance  
tests in open source

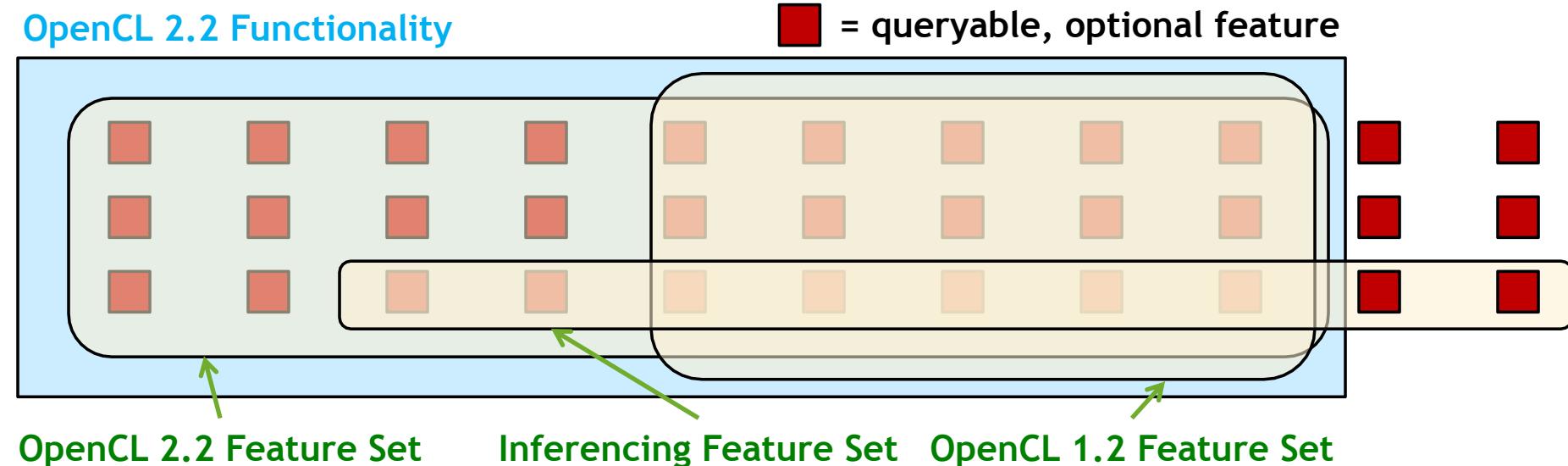
# OpenCL Next

- Enable diverse, heterogenous computing architectures
  - Not just GPUs
- Enable vendors to ship targeted functionality for their customers/markets
  - While still being formally conformant
- Empower middleware + libraries to emulate features
  - Reduces implementation cost
  - Enables portability without vendors implementing all features
- More features in OpenCL become optional - e.g. floating point
  - Enable smaller devices to support recent OpenCL versions
  - Dependency tree between features
- Enhanced query mechanisms
  - Application queries which OpenCL features are supported by a device



# OpenCL Next Feature Sets

- Optional Feature Sets to reduce market fragmentation
  - Opportunity for industry to flexibly agree on commonly supported functionality
  - Features optional in both language and API
- Vendor can support ANY combination of features
  - Not necessary to support any Feature Set - as long as supported features are conformant
  - If features include 100% of a Feature Set - can claim conformance to that Feature Set
  - Supporting popular Feature Sets may help drive sales
  - An Implementation may support multiple Feature Sets



# OpenCL Next and DSPs

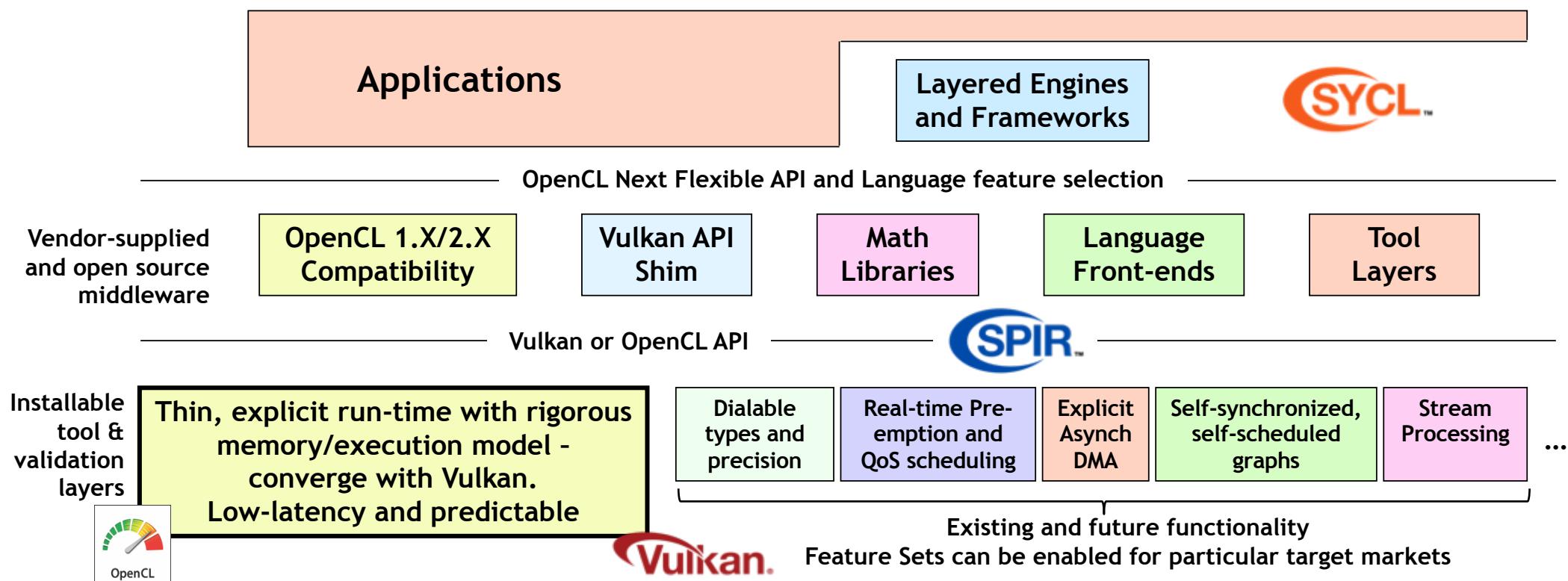
- Flexible Feature Sets enable precise market targeting
  - Embedded imaging, vision and inferencing etc.
- What functionality should be optional?
  - Flexible reduced precision
  - Conformance without IEEE 32 Floating Point
- What new features for DSP and embedded devices?
  - E.g. Explicit synch DMA
  - What else?

OpenCL can be a low-cost, flexible run-time framework to enable cost-effective code generation that is portable code across a wide range of heterogenous devices

	DSP
8-bit int	✓
16-bit int	✓
32-bit int	✓
64-bit int	✗
16-bit float	✗
32-bit float	✗
64-bit float	✗

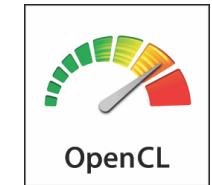
# Longer Term Roadmap Convergence

- Use OpenCL Next flexibility to enable OpenCL API/language over diverse run-times
  - OpenCL, Vulkan or third-party
- Layered ecosystem for backwards-compatibility and market flexibility
  - Feature sets for target market agility



# Comments and Questions?

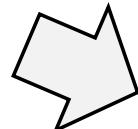
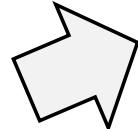
- We value your honest feedback!
- Looking for input/feedback, especially from DSP vendors!



OpenCL 2.2

OpenCL Next defines optional Features and Feature Sets for deployment flexibility

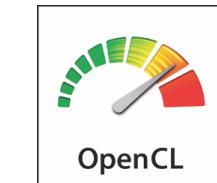
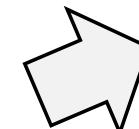
For Heterogenous Compute and Embedded Processors



Help expand core Vulkan Compute - and expose through subsetted OpenCL languages and API



Explore further kernel, language and API convergence



Use Feature and Feature Sets to reduce industry fragmentation and enable new classes of devices to be compliant with the latest OpenCL