# KHRONOS® GROUP
# Fast Forward

**Neil Trevett**
**Khronos President**
**NVIDIA VP Developer Ecosystems**
ntrevett@nvidia.com | @neilt3d
**November 2019**

OpenCL

OpenXR™

Vulkan®

SYCL™

glTF™

SPIR™

OpenVX™

NNEF™

SIGGRAPH ASIA 2019 BRISBANE

# What are Open Standards?

Interoperability Standards define an agreed communication protocol between two 'entities'

Common products use 100s of open standards



Device to Wireless Networks

Device to its Charger

Internal Components to other Internal Components

Downloaded Web content to the Web Browser

Camera App to Video and Photo Playback Applications

Games to 3D GPU Acceleration

## Many Standard Defining Organizations (SDOs)
Each has a focus area of expertise that gathers an effective quorum
Each creates a safe space for cooperation
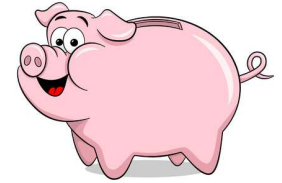
# The Need for Interoperability Standards

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

**Standards Reduce Costs**
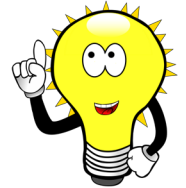By sharing development between many companies and driving volume

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

**Standards Do Not Stifle Innovation**
Companies can compete on implementation quality, performance, power etc. etc.

**True OPEN Standards**
Are not controlled by a single company – but by the industry – typically through an SDO
Well defined participation, governance and intellectual property frameworks

**KHRONOS** GROUP

>150 Members ~ 40% US, 30% Europe, 30% Asia

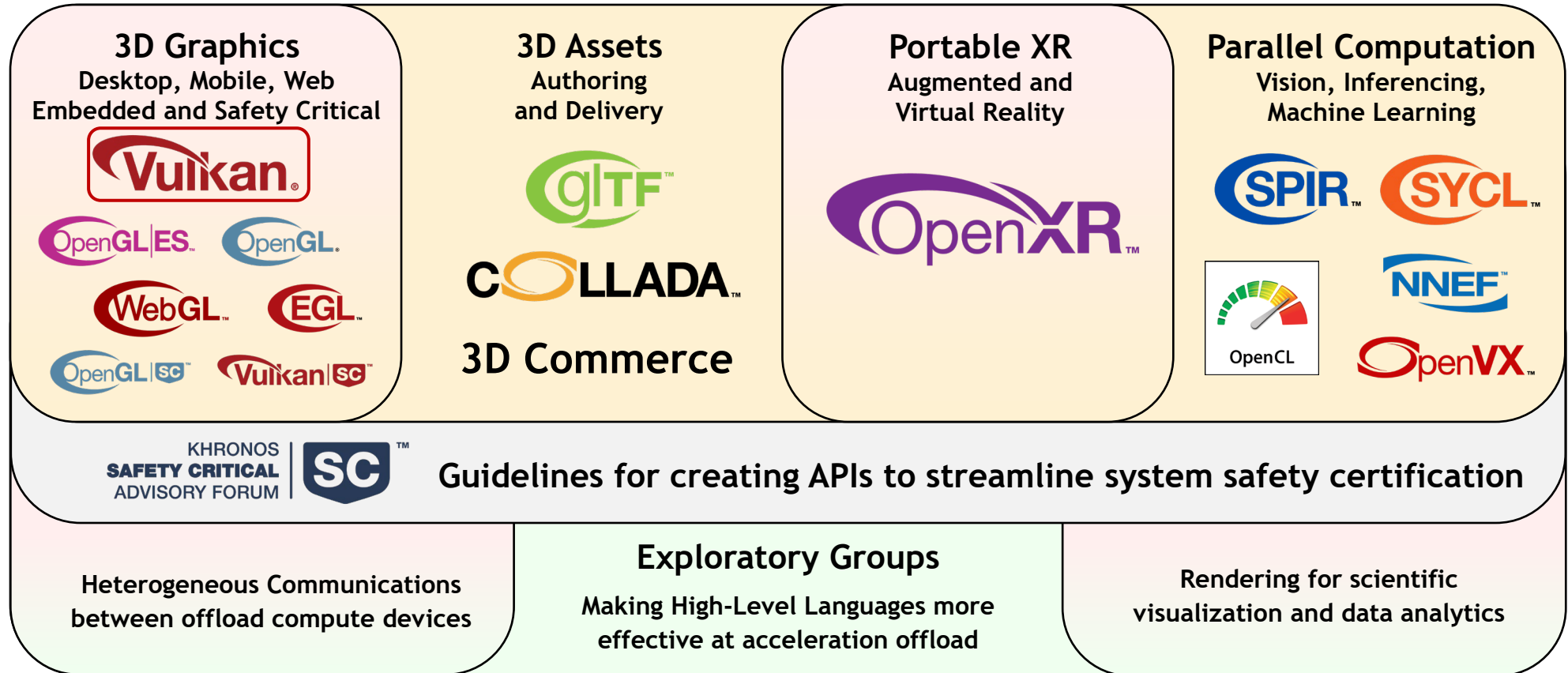Khronos is an open, non-profit, member-driven industry consortium developing royalty-free standards to harness the power of silicon acceleration for demanding graphics rendering and computationally intensive applications such as 3D Graphics, Virtual Reality, Augmented Reality, and Machine Learning

# Khronos Asia Pacific Members



**Khronos warmly welcomes Australian and Asian company participation!!**

# Khronos Active Initiatives

**3D Graphics**
Desktop, Mobile, Web
Embedded and Safety Critical

Vulkan.
OpenGL|ES OpenGL.
WebGL EGL
OpenGL|SC Vulkan|SC

**3D Assets**
Authoring
and Delivery

glTF
COLLADA.

**3D Commerce**

**Portable XR**
Augmented and
Virtual Reality

OpenXR.

**Parallel Computation**
Vision, Inferencing,
Machine Learning

SPIR SYCL
NNEF
OpenCL OpenVX

KHRONOS
SAFETY CRITICAL **SC**
ADVISORY FORUM
Guidelines for creating APIs to streamline system safety certification

**Heterogeneous Communications**
between offload compute devices

**Exploratory Groups**
Making High-Level Languages more
effective at acceleration offload

**Rendering for scientific**
visualization and data analytics

# Vulkan for Direct GPU Control

**Complex drivers cause overhead and inconsistent behavior across vendors**

**Always active error handling**

**Full GLSL preprocessor and compiler in driver**

**OpenGL vs. OpenGL ES**

OpenGL|ES
OpenGL

**Application**
Single thread per context

**High-level Driver Abstraction**
**Layered GPU Control**
Context management
Memory allocation
Full GLSL compiler
Error detection

**GPU**

A Graphics API

Vulkan

**Application**
Memory allocation
Thread management
Explicit Synchronization
Multi-threaded generation of command buffers

**Multiple Front-end Compilers**
GLSL, HLSL etc.

SPIR
SPIR-V
pre-compiled shaders

**Thin Driver**
**Explicit GPU Control**

**Loadable debug and validation layers**

**GPU**

A GPU API

**Simpler drivers - application has the best knowledge for holistic optimization – no 'driver magic'**

**Explicit creation of API objects before usage – efficient, predictable execution**

**Easier portability - no fighting with different vendor heuristics**

**Validation and debug layers loaded only when needed**

**SPIR-V intermediate language: shading language flexibility**

**Unified API across mobile and desktop platforms**

**Multiple graphics, command and DMA queues**
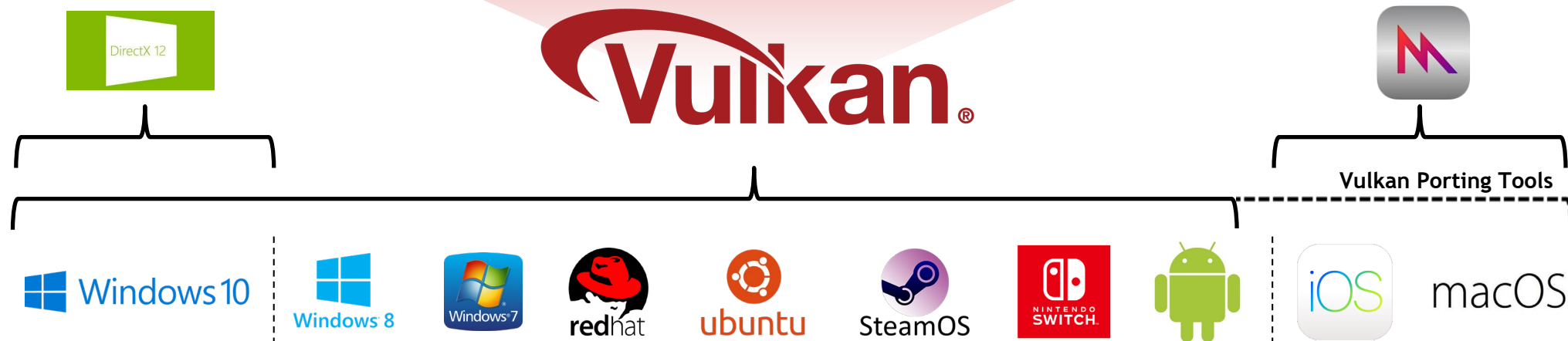
# Vulkan and New Generation GPU APIs

Modern architecture | Low overhead| Multi-thread friendly
EXPLICIT GPU access for EFFICIENT, LOW-LATENCY,
PREDICTABLE performance

**Vulkan Porting Tools**

**Vulkan is a non-proprietary, royalty-free open standard**
**Portable across multiple platforms – desktop, mobile and embedded**

# Pervasive Vulkan

## Major GPU Companies supporting Vulkan for Desktop and Mobile Platforms

AMD · arm · BROADCOM · Imagination · intel · NVIDIA · QUALCOMM · VeriSilicon

## Platforms

**Desktop**
Windows 10 / Windows 8 / Windows 7 / Linux

**Android**
(Android 7.0+)
(Vulkan 1.1 required on Android Q)

**Apple**
(via porting layers)
iOS macOS

**Media Players**

**Consoles**
NINTENDO SWITCH

**Virtual Reality**
Windows 8 / Windows 10 / Windows 7 / Linux

**Cloud Services**
Windows 10 / Linux

**Game Streaming**
NVIDIA GEFORCE NOW / STADIA

**Embedded**

## Game Engines

EPIC GAMES · id · CRYENGINE · unity · VALVE · Croteam Serious Engine · XENKO · NetEase Games Passion of gamers
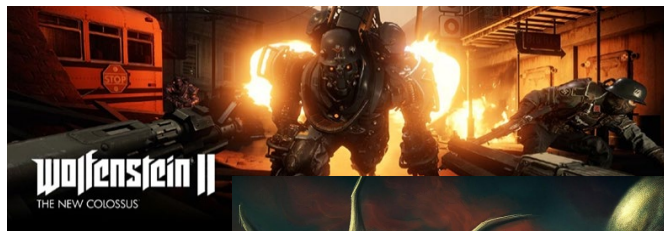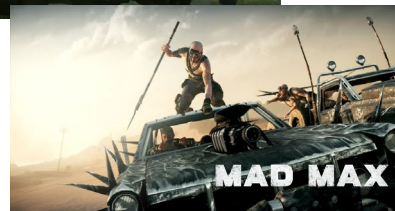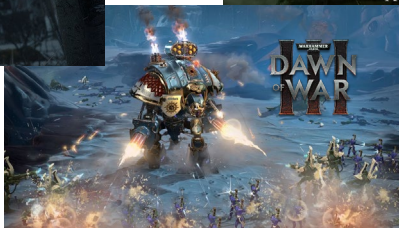
# Vulkan AAA Content Shipping on Desktop…



Vulkan-only AAA Titles on PC

AAA titles on Linux

Titles on PC AND macOS

# ...and Mobile

GALAXY ON FIRE 3

FERAL
MAKE YOUR PLAY

GRID

FORTNITE

FORTNITE
BATTLE ROYALE

EPIC GAMES

Vulkan

Plus....
Lineage 2 Revolution
Heroes of Incredible Tales
Dream League Soccer...

FLARE GAMES

COATSINK
netmarble Games

SUPER EVIL MEGACORP

CODEMASTERS

DIGITAL LEGENDS ENTERTAINMENT

NEXON

FIRST TOUCH GAMES

# Vulkan 1.1 Ecosystem Evolution

## Strengthening Tools and Compilers

Improved developer tools (SDK, validation/debug layers)
Shader toolchain improvements (size, speed, robustness)
Shading language flexibility – HLSL and OpenCL C support
More rigorous conformance testing

**February 2016**
**Vulkan 1.0**

### Vulkan 1.0 Extensions
Maintenance updates plus additional functionality

Multiview
Multi-GPU
Enhanced Windows System Integration
Increased Shader Flexibility:
16-bit storage, Variable Pointers
Enhanced Cross-Process and
Cross-API Sharing

**March 2018**
**Vulkan 1.1**

**Integration of 1.0 Extensions
plus new functionality
e.g. Subgroup Operations**

## Building Vulkan's Future

Listen and prioritize developer needs
Drive GPU technology

### Released Vulkan 1.1 Extensions
Reduced precision arithmetic types in shaders
Bindless resources
HLSL-compatible memory layouts
Formal memory model
Buffer references
Timeline semaphores
OpenGL-class lines and Interop
https://www.khronos.org/registry/vulkan/specs/1.1-khr-extensions/html/vkspec.html#extension-appendices-list

### Roadmap Discussions
Machine Learning acceleration
Ray Tracing
Video encode / decode
Generalized subgroup operations

## Widening Platform Support

Pervasive GPU vendor native driver availability
Open source drivers – ANV (Intel), AMDVLK/RADV (AMD)
Vulkan Portability to macOS/iOS and DX12

# OpenGL Vulkan Interop

- **Enables OpenGL applications to incrementally leverage Vulkan functionality**
  - Shared explicit memory objects

- **Dassault Systèmes achieves interactive object space AO in CATIA, an OpenGL application**
  - Using the NVIDIA Vulkan VKRay vendor extension for Ray Tracing
  - See the Demo at the NVIDIA booth

# Key Vulkan Online Open Source Resources

**Vulkan Samples**
Collection of samples and resources to aid developing optimized Vulkan applications
https://github.com/KhronosGroup/Vulkan-Samples

**Vulkan Guide**
Help for developers to get up and going with the world of Vulkan with kinks to many other useful resources
https://github.com/KhronosGroup/Vulkan-Guide

**RenderDoc Debugger**
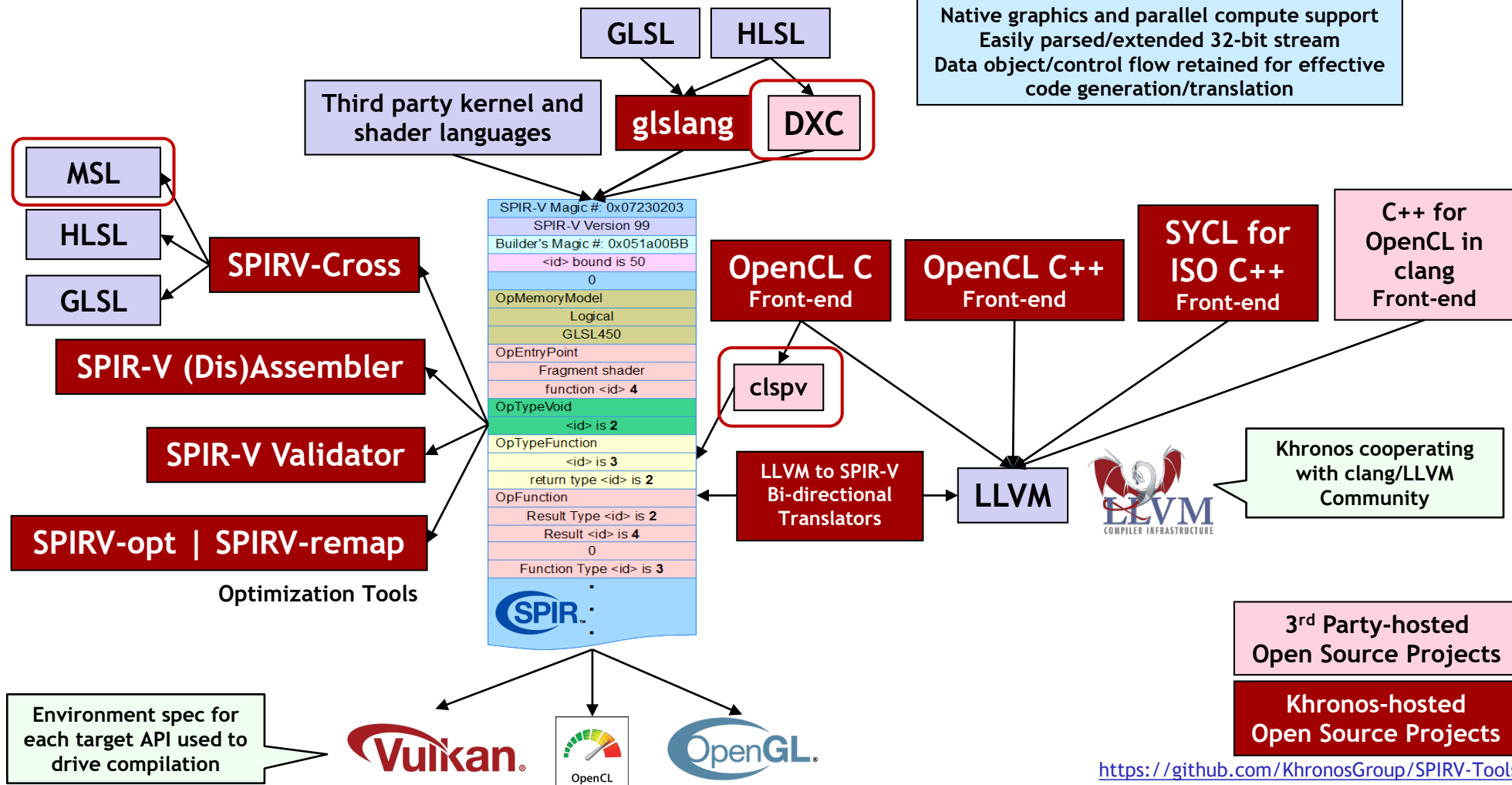Single-frame capture and detailed introspection of any application
https://renderdoc.org/

**Vulkan SDK with Development/Debug Layers**
Windows, Linux - Ubuntu packages, Linux- Tarball, macOS
www.vulkan.lunarg.com

# SPIR-V Ecosystem

GLSL  HLSL

**SPIR-V**
Khronos-defined cross-API IR
Native graphics and parallel compute support
Easily parsed/extended 32-bit stream
Data object/control flow retained for effective
code generation/translation

Third party kernel and shader languages

**glslang**  **DXC**

MSL

HLSL

**SPIRV-Cross**

GLSL

**SPIR-V (Dis)Assembler**

**SPIR-V Validator**

**SPIRV-opt | SPIRV-remap**

Optimization Tools

| SPIR-V Magic #: 0x07230203 |
| SPIR-V Version 99 |
| Builder's Magic #: 0x051a00BB |
| <id> bound is 50 |
| 0 |
| OpMemoryModel |
| Logical |
| GLSL450 |
| OpEntryPoint |
| Fragment shader |
| function <id> **4** |
| OpTypeVoid |
| <id> is **2** |
| OpTypeFunction |
| <id> is **3** |
| return type <id> is **2** |
| OpFunction |
| Result Type <id> is **2** |
| Result <id> is **4** |
| 0 |
| Function Type <id> is **3** |

**OpenCL C**
Front-end

**OpenCL C++**
Front-end

**SYCL for
ISO C++**
Front-end

C++ for
OpenCL in
clang
Front-end

**clspv**

**LLVM to SPIR-V
Bi-directional
Translators**

**LLVM**

Khronos cooperating
with clang/LLVM
Community

LLVM
COMPILER INFRASTRUCTURE

**3rd Party-hosted
Open Source Projects**

**Khronos-hosted
Open Source Projects**

Environment spec for
each target API used to
drive compilation

Vulkan.  OpenCL  OpenGL.

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

# Open Source Layering Projects

**Breaking through platform fragmentation**

Vulkan added OpenGL-style line extension

Vulkan adding more compute for fuller support for OpenCL kernel deployment

Vulkan added extensions to ease layering of DX

Vulkan is an effective porting layer for app portability and stack simplification

| Layers Over | Vulkan | OpenGL | OpenCL | OpenGL ES | DX12 | DX9-11 |
|---|---|---|---|---|---|---|
| **Vulkan** | | Zink | clspv clvk | GLOVE Angle | vkd3d | DXVK D9VK |
| **OpenGL** | gfx-rs Ashes | | | Angle | | |
| **DX12** | gfx-rs | | | | | |
| **DX9-11** | gfx-rs Ashes | | | Angle | | |
| **Metal** | MoltenVK gfx-rs | | | Angle | | |

Vulkan Portability enables multi-vendor layered subsets to be queryable and conformant

Growing interest to offset Apple deprecation?

Demand for Vulkan everywhere, even if no native drivers on platform

Need for consistent OpenGL ES everywhere, primarily for WebGL

# Vulkan Portability Initiative

**Enabling Vulkan applications on platforms without native drivers by layering cleanly queryable subsets of Vulkan over DX12, Metal and other APIs**

## Multiple Layered Vulkan Implementations

Additional open source run-times over additional backends
E.g. gfx-rs for Vulkan over Metal and DX12 - useful for Vulkan on UWP platforms such as Windows 10 S, Polaris, Xbox One.
Secondary backends include OpenGL/D3D11
https://github.com/gfx-rs/gfx
https://github.com/gfx-rs/portability

## Portability Extension

Layered implementations can portably expose
what Vulkan functionality is not supported

**TODAY**

Open source tools, SDKs and libraries to bring Vulkan 1.0 applications to Apple using Metal

## Extend Vulkan Conformance Test Suite

To handle layered implementations – what is present must work!

## Enhanced Vulkan Layers

Extend DevSim/Validation Layers to flag or simulate queries for features not present

# Vulkan Portability Initiative on Apple

Almost all mandatory Vulkan 1.0 functionality is supported:
No Triangle Fans
No separate stencil reference masks

Selected Optional Features and Extensions are added as required - driven by industry input and feedback
Robust buffer access
BC texture compressed formats
Fragment shader atomics
Tesselation
https://github.com/KhronosGroup/MoltenVK

**Applications**

Open source SDK to build, run, and debug applications on macOS - including validation layer support
https://vulkan.lunarg.com/

**Vulkan macOS SDK**

**SPIRV-Cross**
Convert SPIR-V shaders to Metal Shaders

**macOS / iOS Run-time**
Maps Vulkan to Metal

Open source beta release for macOS

MoltenVK supports macOS 10.11 / iOS 9.0 and up

Open source for MacOS and iOS
Free to use - no fees or royalties including commercial apps

# Vulkan Apps Shipping On Apple

**Vulkan** PORTABILITY

**Forsaken Remastered** was just updated with **Vulkan** support! If you're on Linux, you're probably hitting 60fps with the existing OpenGL renderer, but it's good to be future proof. If you're on a Mac, though, you *definitely* want to switch. On my MacBook, the framerate goes from around 15 to a solid 60!

### Initial Vulkan Performance On macOS With Dota 2 Is Looking Very Good
Written by Michael Larabel in Valve on 1 June 2018 at 05:37 PM EDT. 34 Comments

Yesterday Valve released Vulkan support for Dota 2 on macOS. Indeed, this first major game relying upon MoltenVK for mapping Vulkan over the Apple Metal drivers is delivering performance gains.

### Valve Releases Artifact As Its Cross-Platform, Vulkan-Powered Digital Card Game
Written by Michael Larabel in Valve on 28 November 2018 at 04:16 PM EST. 29 Comments

Valve managed to ship their latest game today as planned and without any major delays.

Artifact is now available with launch-day support for Linux, macOS, and Windows. Artifact is a competitive digital card game is targeting Dota 2 players as well as card gaming enthusiasts. Valve still plans to evolve Artifact and its gameplay more

**Production Dota 2 on Mac Ships – up to 50% more perf than Apple's OpenGL**

**First iOS Apps using MoltenVK ship through app store**

**Qt Running on Mac through MoltenVK**

**Multiple iOS and macOS apps shipping e.g. Forsaken Remastered**

**Google Filament PBR Renderer on Mac**

**Initial ports of DX games in progress using Vulkan on Mac**

A R T I F A C T

**Artifact from Steam ships on MoltenVK on macOS - first Vulkan-only Valve app on Mac**

RPCS3

**RPCS3 PlayStation 3 Emulator on Mac**

**Dolphin GameCube and Wii Emulator working on MacOS**

**Diligent Engine runs on MacOS**

DOTA UNDERLORDS

**Artifact from Steam ships on MoltenVK on macOS - second Vulkan-only Valve app on Mac**

| June 2018 | September 2018 | November 2018 | January 2019 | June 2019 |

# Running DX Games on Linux Over Vulkan

- **DXVK – Direct3D 10/11 emulator running over Vulkan**
  - Open source on GitHub - developed by Philip Rebohle with support from Valve

- **Vulkan has added multiple extensions to support efficient layering of D3D**
  - Removing impedance mismatches between the two APIs

- **DXVK, Wine Windows Compatibility Layer and Valve Proton tool**
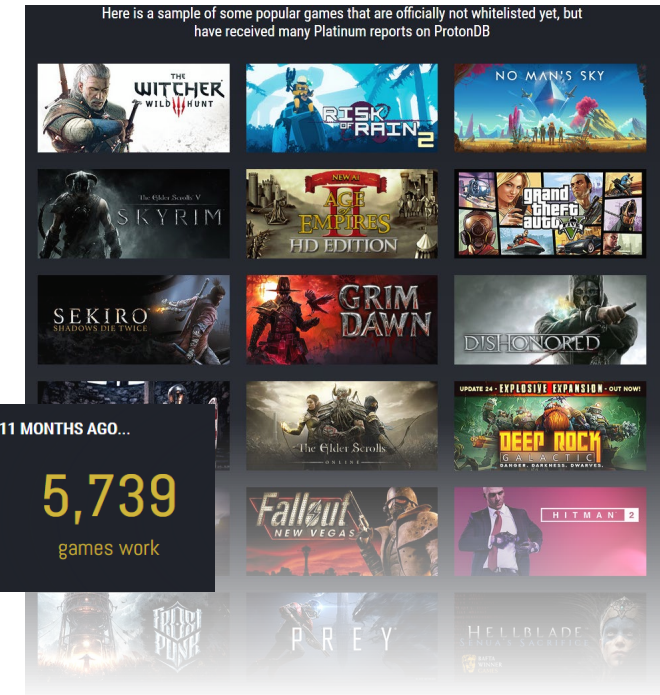  - Enable thousands of PC games on Linux

**Extensions created in response to DXVK issues**
VK_EXT_transform_feedback
VK_EXT_depth_clip_enable
VK_EXT_host_query_reset
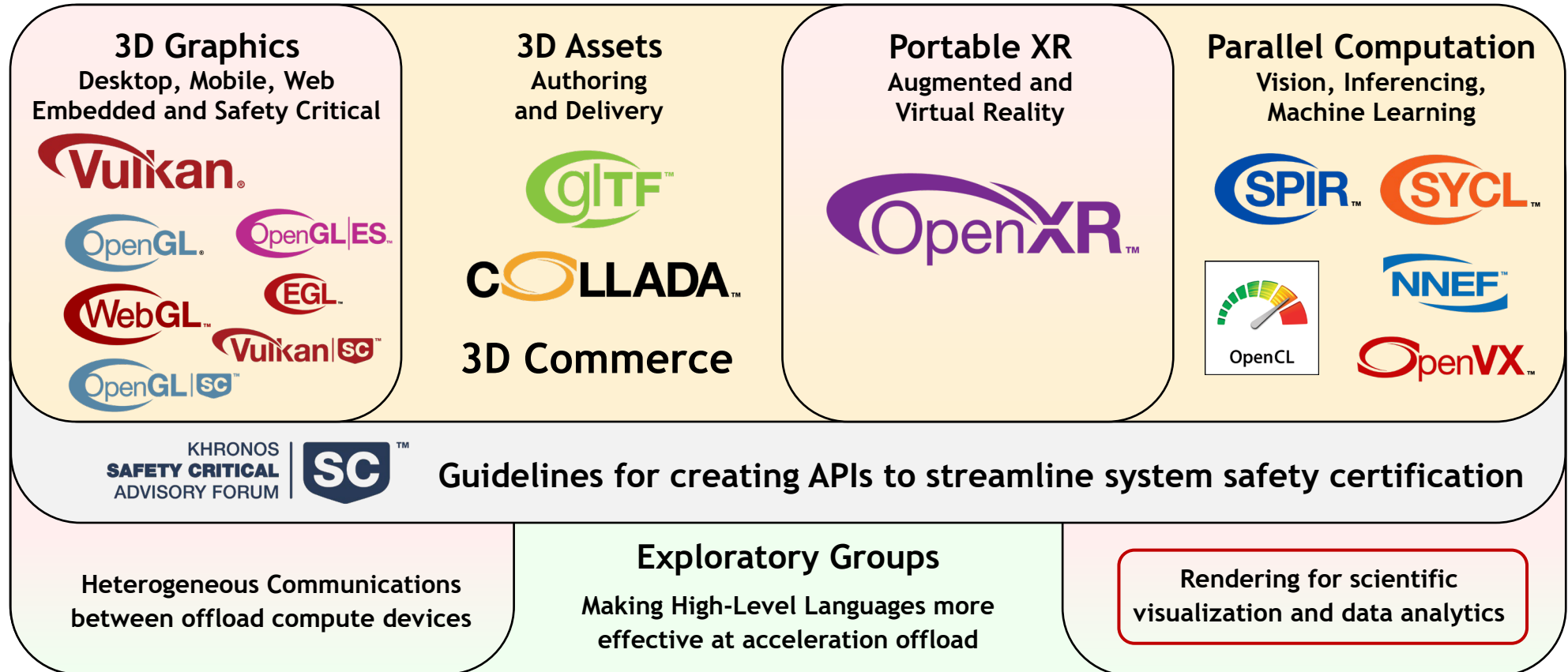VK_EXT_texel_buffer_alignment
VK_EXT_shader_demote_to_helper_invocation

**Other extensions used by DXVK**
VK_EXT_conditional_rendering
VK_EXT_memory_budget
VK_EXT_memory_priority
VK_EXT_shader_viewport_index_layer
VK_EXT_vertex_attribute_divisor
VK_KHR_draw_indirect_count
VK_KHR_shader_draw_parameters

https://www.protondb.com

Here is a sample of some popular games that are officially not whitelisted yet, but have received many Platinum reports on ProtonDB

SINCE THE RELEASE OF PROTON ONLY **11 MONTHS AGO**...

**51,447** reports written  **8,833** games reported  **5,739** games work

# Khronos Active Initiatives

**3D Graphics**
Desktop, Mobile, Web
Embedded and Safety Critical

**Vulkan**
**OpenGL** **OpenGL ES**
**WebGL** **EGL**
**Vulkan SC**
**OpenGL SC**

**3D Assets**
Authoring
and Delivery

**glTF**
**COLLADA**

**3D Commerce**

**Portable XR**
Augmented and
Virtual Reality

**OpenXR**

**Parallel Computation**
Vision, Inferencing,
Machine Learning

**SPIR** **SYCL**
**OpenCL** **NNEF**
**OpenVX**

**KHRONOS SAFETY CRITICAL ADVISORY FORUM** | **SC**   Guidelines for creating APIs to streamline system safety certification

Heterogeneous Communications
between offload compute devices

**Exploratory Groups**
Making High-Level Languages more
effective at acceleration offload

Rendering for scientific
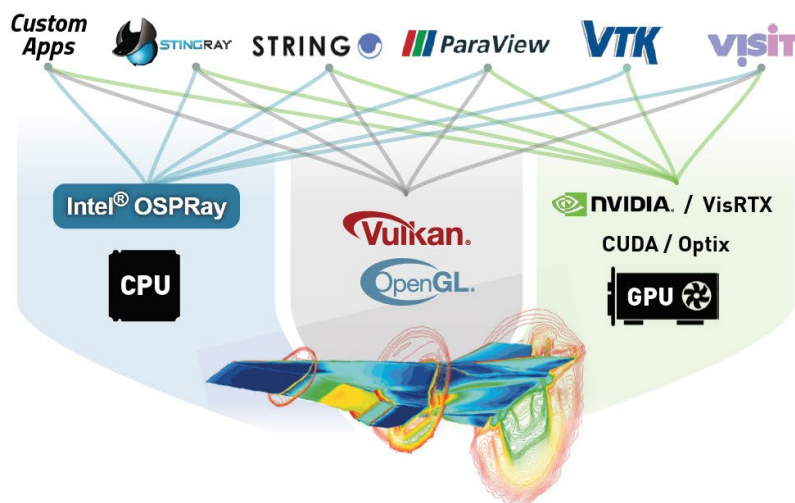visualization and data analytics

# Analytic Rendering Exploratory Group

**Analytic Rendering is image generation performed primarily to gain and communicate insights into complex data sets primarily for scientific visualization and data analytics**
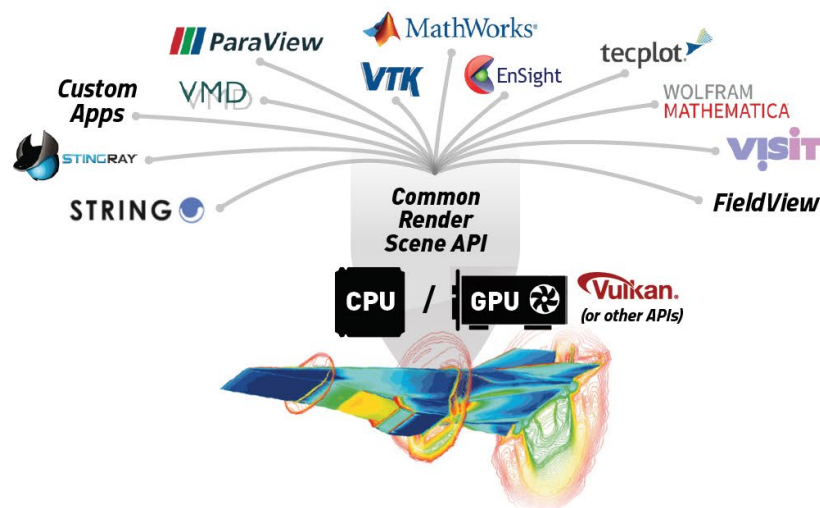
**Is there a need for a cross-platform open standard API?**



**Visualization Apps and Engines have to be ported to multiple APIs**

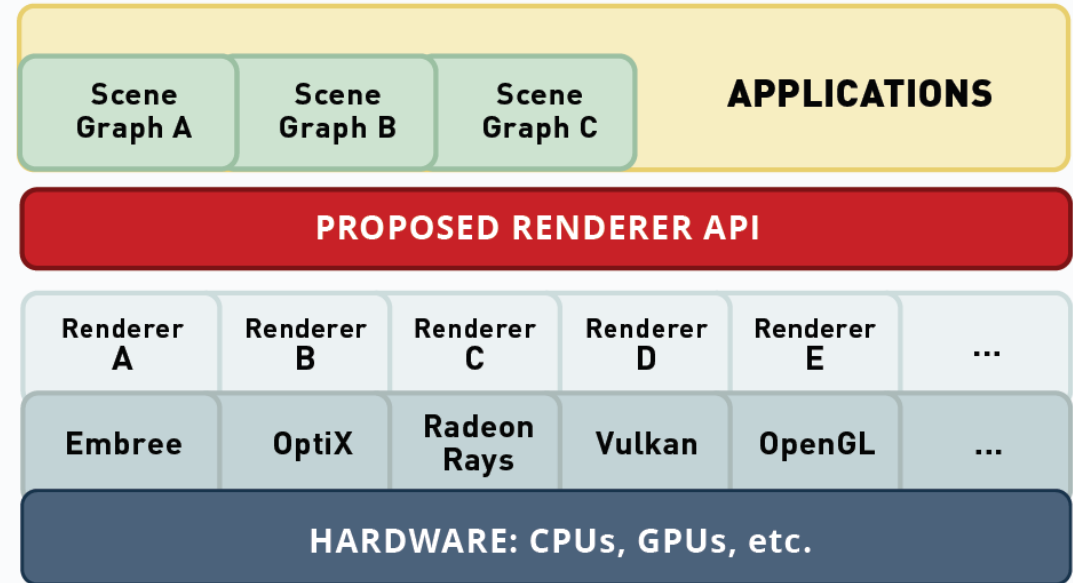**Cross-vendor API to provide access to state-of-the-art rendering across multiple platforms**
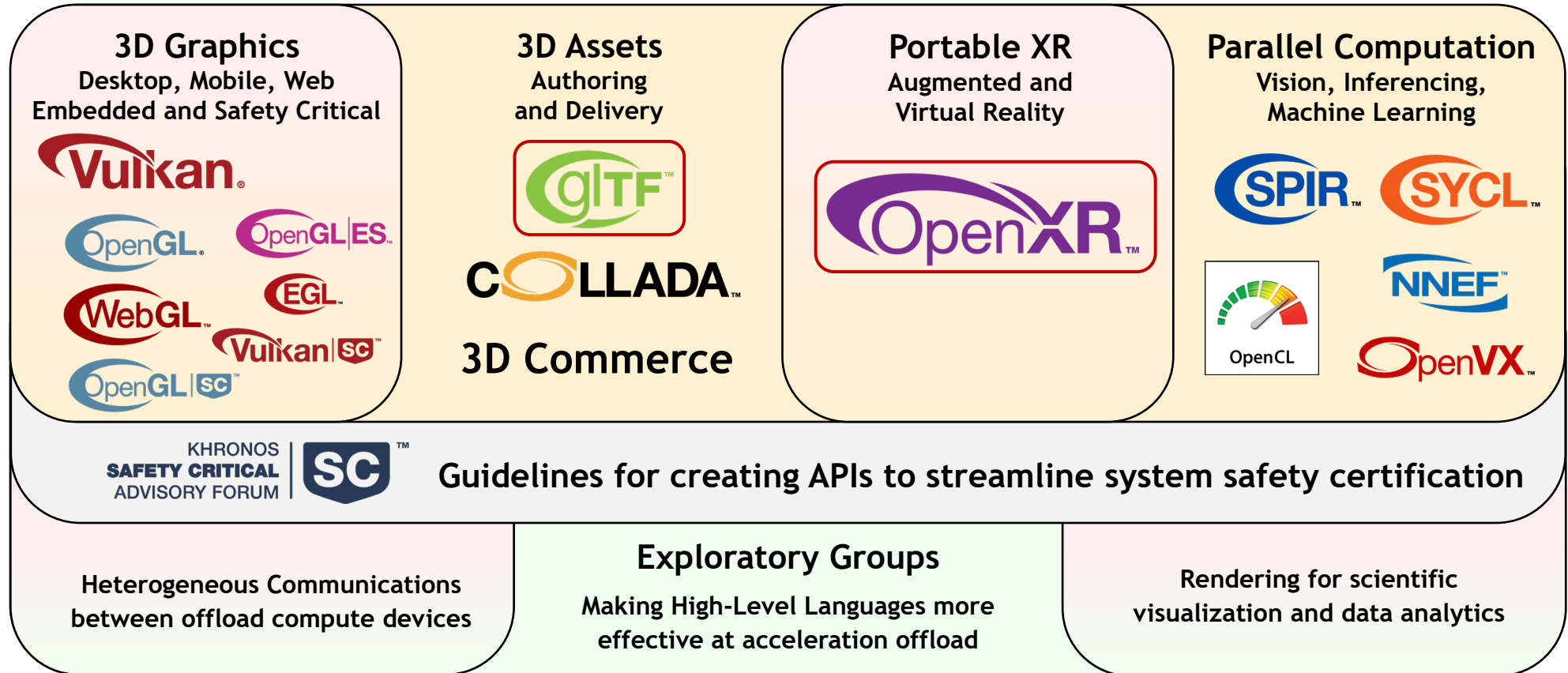
# Potential Analytic Rendering API Design

Rather than specifying the details of the rendering process, an Analytic Rendering API would enable a visualization application to simply describe the relationship between objects in a scene to be rendered and leave the details of the rendering process to a backend renderer

| Scene Graph A | Scene Graph B | Scene Graph C | **APPLICATIONS** |
|---|---|---|---|

**PROPOSED RENDERER API**

| Renderer A | Renderer B | Renderer C | Renderer D | Renderer E | ... |
|---|---|---|---|---|---|
| Embree | OptiX | Radeon Rays | Vulkan | OpenGL | ... |

**HARDWARE: CPUs, GPUs, etc.**

delta h
Ingenieurgesellschaft
(intel)
Kitware
OAK RIDGE
National Laboratory

NVIDIA
SURVICE
ENGINEERING COMPANY
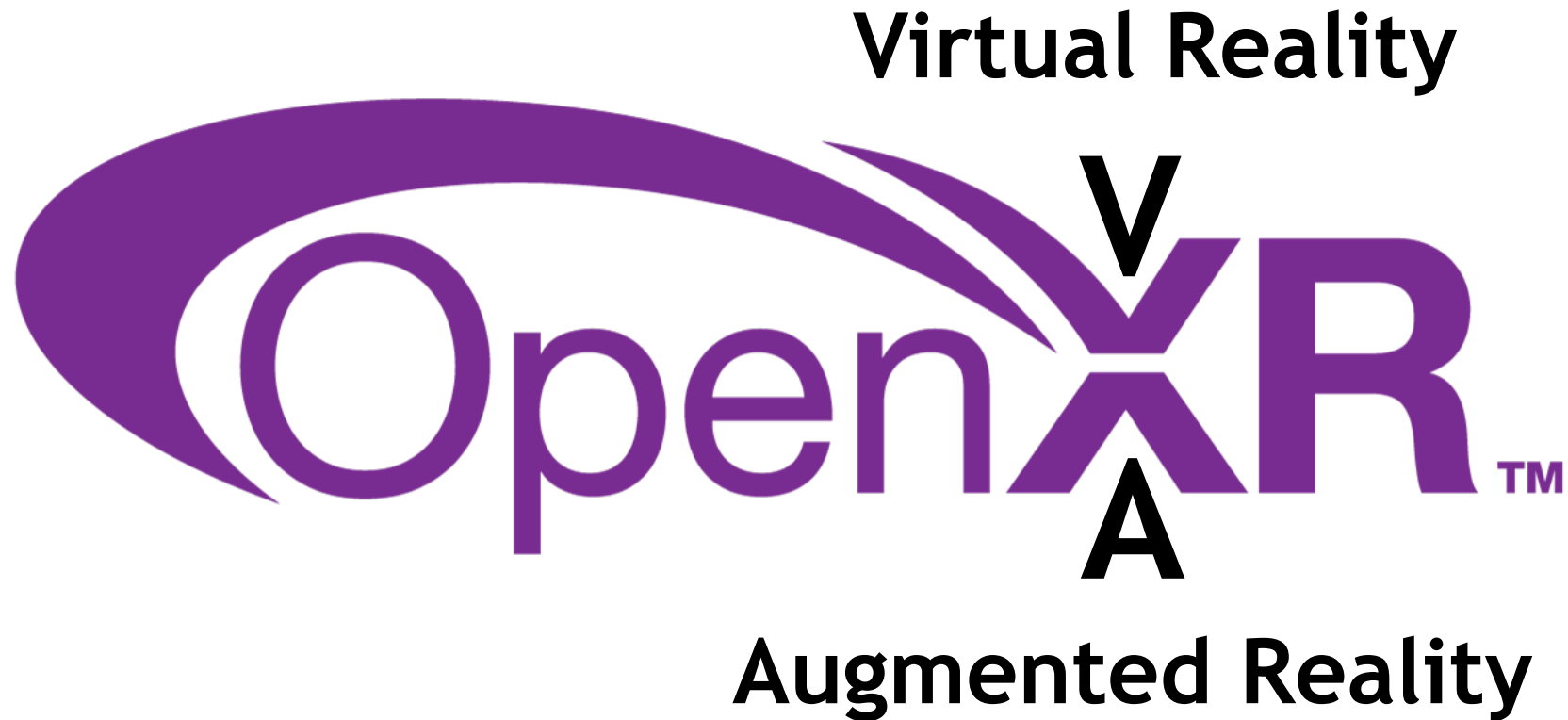TACC

**Some Initial Exploratory Group Members**

**Khronos Exploratory Groups discuss the need for a new standard with no cost or IP Implications**
**Open to all – even non-members - more details**
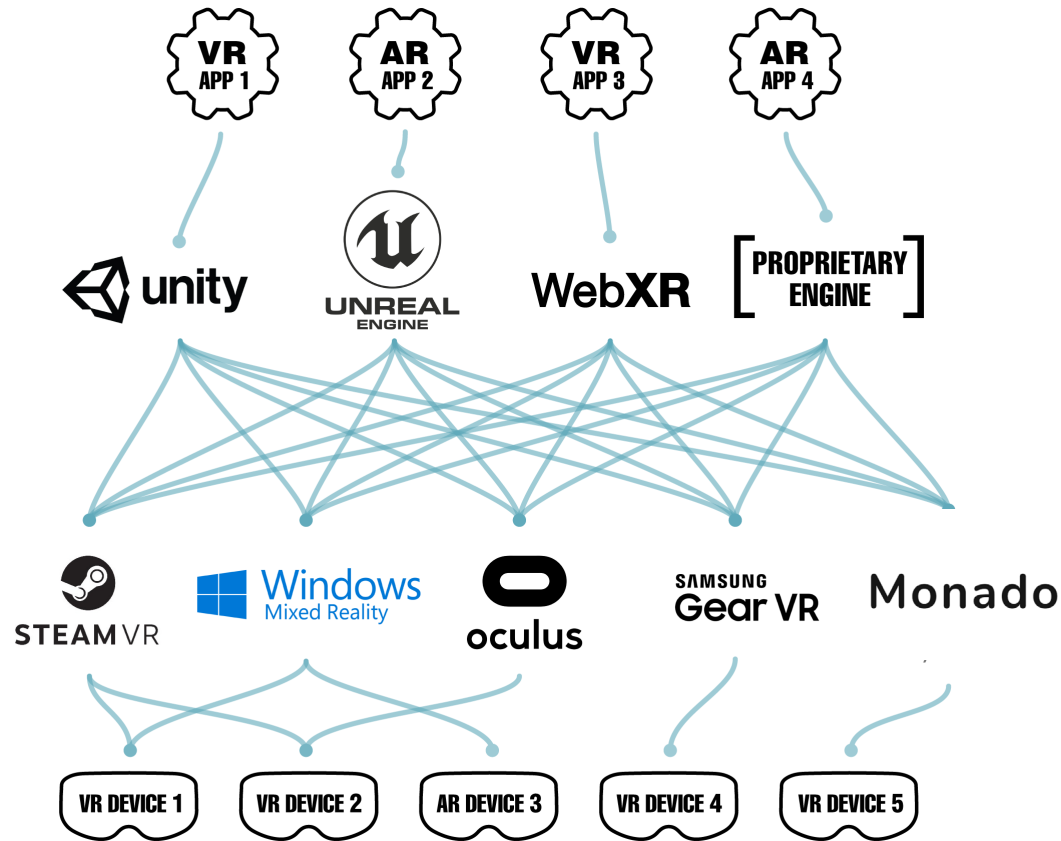https://www.khronos.org/exploratory/analytic-rendering/

# Khronos Active Initiatives

**3D Graphics**
Desktop, Mobile, Web
Embedded and Safety Critical

**Vulkan**
**OpenGL** **OpenGL|ES**
**WebGL** **EGL**
**Vulkan|SC**
**OpenGL|SC**

**3D Assets**
Authoring
and Delivery

**glTF**
**COLLADA**

**3D Commerce**

**Portable XR**
Augmented and
Virtual Reality

**OpenXR**

**Parallel Computation**
Vision, Inferencing,
Machine Learning

**SPIR** **SYCL**
**OpenCL** **NNEF**
**OpenVX**

**KHRONOS SAFETY CRITICAL ADVISORY FORUM | SC** Guidelines for creating APIs to streamline system safety certification

**Heterogeneous Communications**
between offload compute devices

**Exploratory Groups**
Making High-Level Languages more
effective at acceleration offload

Rendering for scientific
visualization and data analytics

# XR = AR + VR

OpenXR provides cross-platform, high-performance
access to AR and VR platforms and devices

**Virtual Reality**
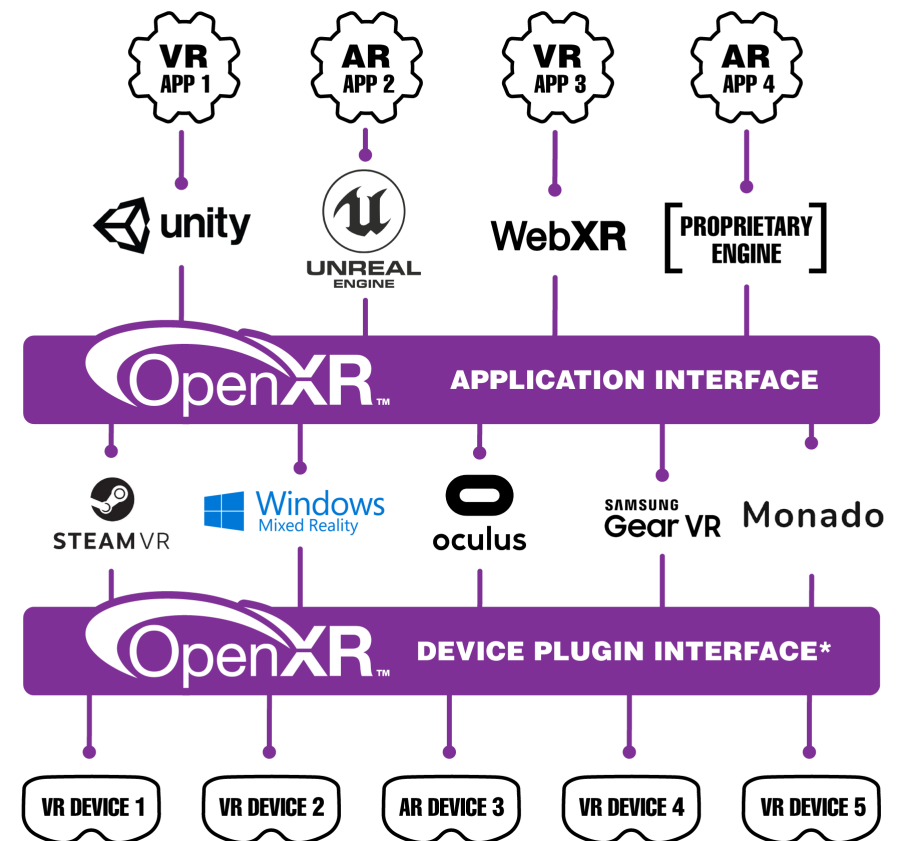
OpenXR™

**Augmented Reality**

# OpenXR – Solving XR Fragmentation



**Before OpenXR**
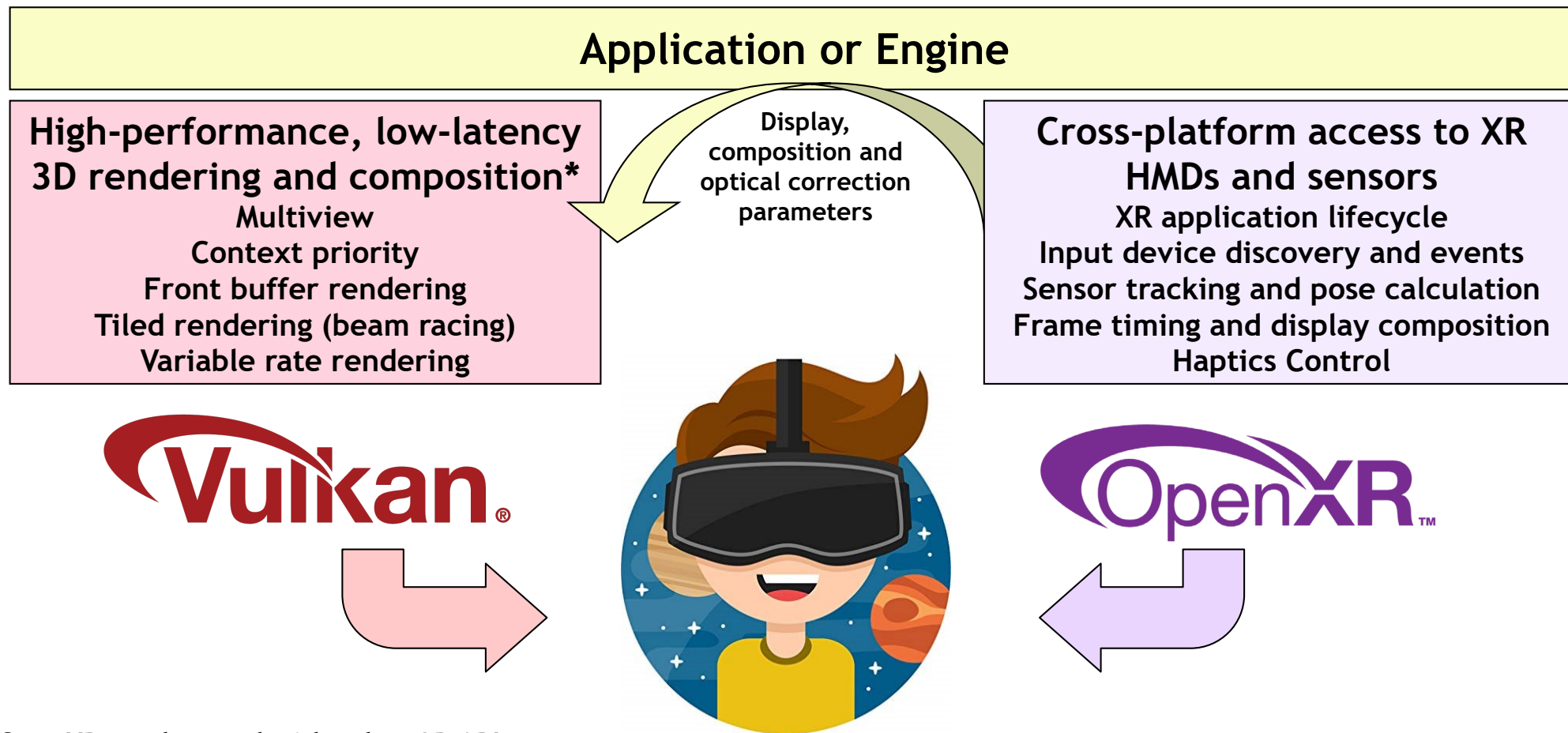
XR Market Fragmentation

**After OpenXR**

Wide interoperability of XR apps and devices

* OpenXR 1.0 is focused on enabling cross-platform applications. Optional device plugin interface will be supported post V1.0
** Check OpenXR Landing Page for exact availabiliy of OpenXR in shipping run-times and devices www.khronos.org/openxr

# OpenXR is used with a 3D API

**Application or Engine**

**High-performance, low-latency 3D rendering and composition***
Multiview
Context priority
Front buffer rendering
Tiled rendering (beam racing)
Variable rate rendering

Display, composition and optical correction parameters

**Cross-platform access to XR HMDs and sensors**
XR application lifecycle
Input device discovery and events
Sensor tracking and pose calculation
Frame timing and display composition
Haptics Control

**Vulkan®**

**OpenXR™**

**\* OpenXR can be used with other 3D APIs such as Direct3D, OpenGL and OpenGL ES**

# Companies Publicly Supporting OpenXR

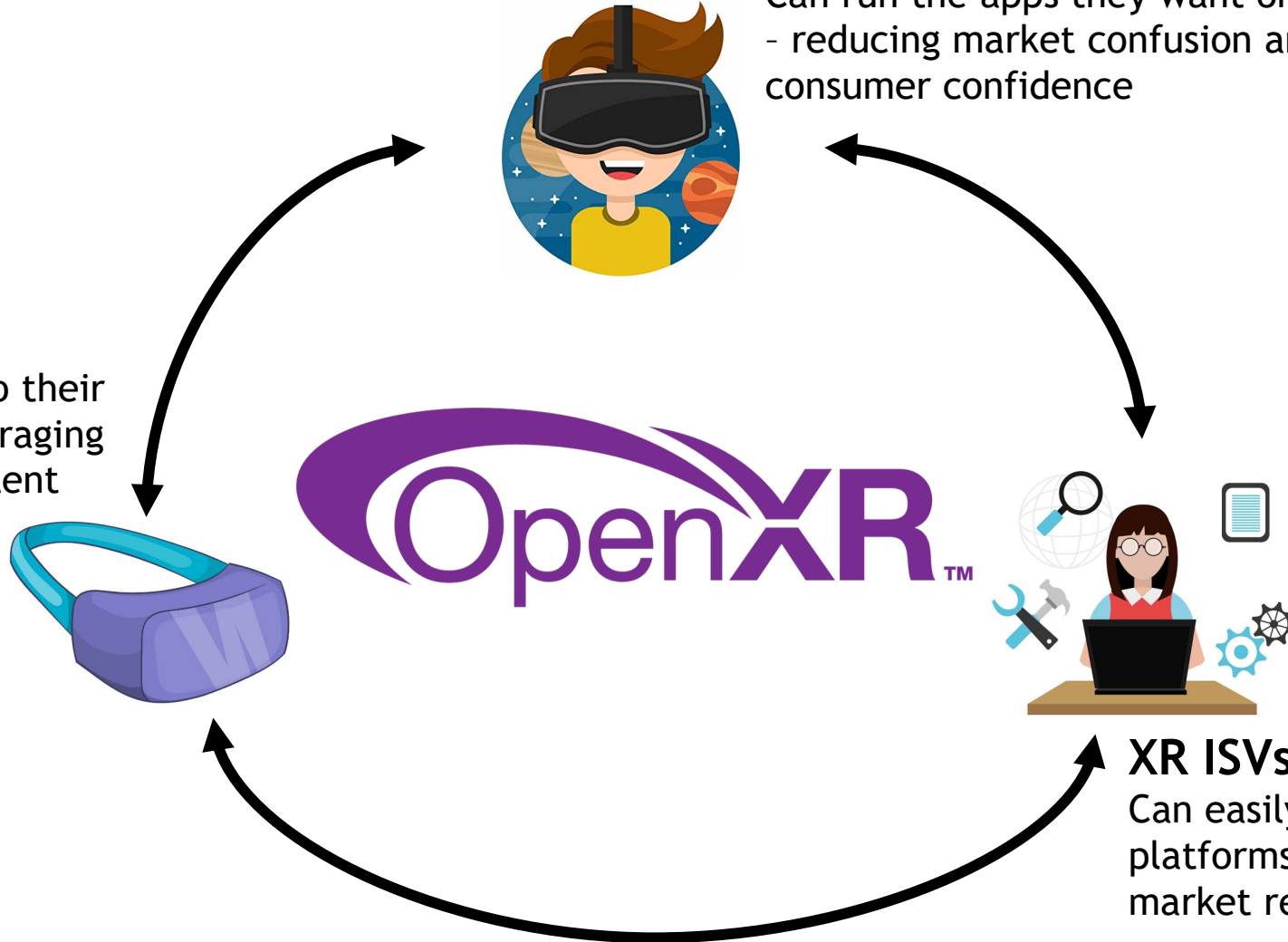AMD · antilatency · arm · AREA · C-O COLLABORA · CTRL-labs · DisplayLink XR

EPIC GAMES · HTC · Google · hp · HUAWEI · Imagination · intel · LG · logitech

LUNAR G · magic leap · MEDIATEK · Microsoft · mozilla · NOKIA · NVIDIA

oculus · Pico · pluto · QUALCOMM · RAZER · SAMSUNG · SONY

tobii · unity · VALVE · VARJO · VeriSilicon · 兆芯 · zSpace

**OpenXR is a collaborative design**
**Integrating many lessons from proprietary 'first-generation' XR API designs**
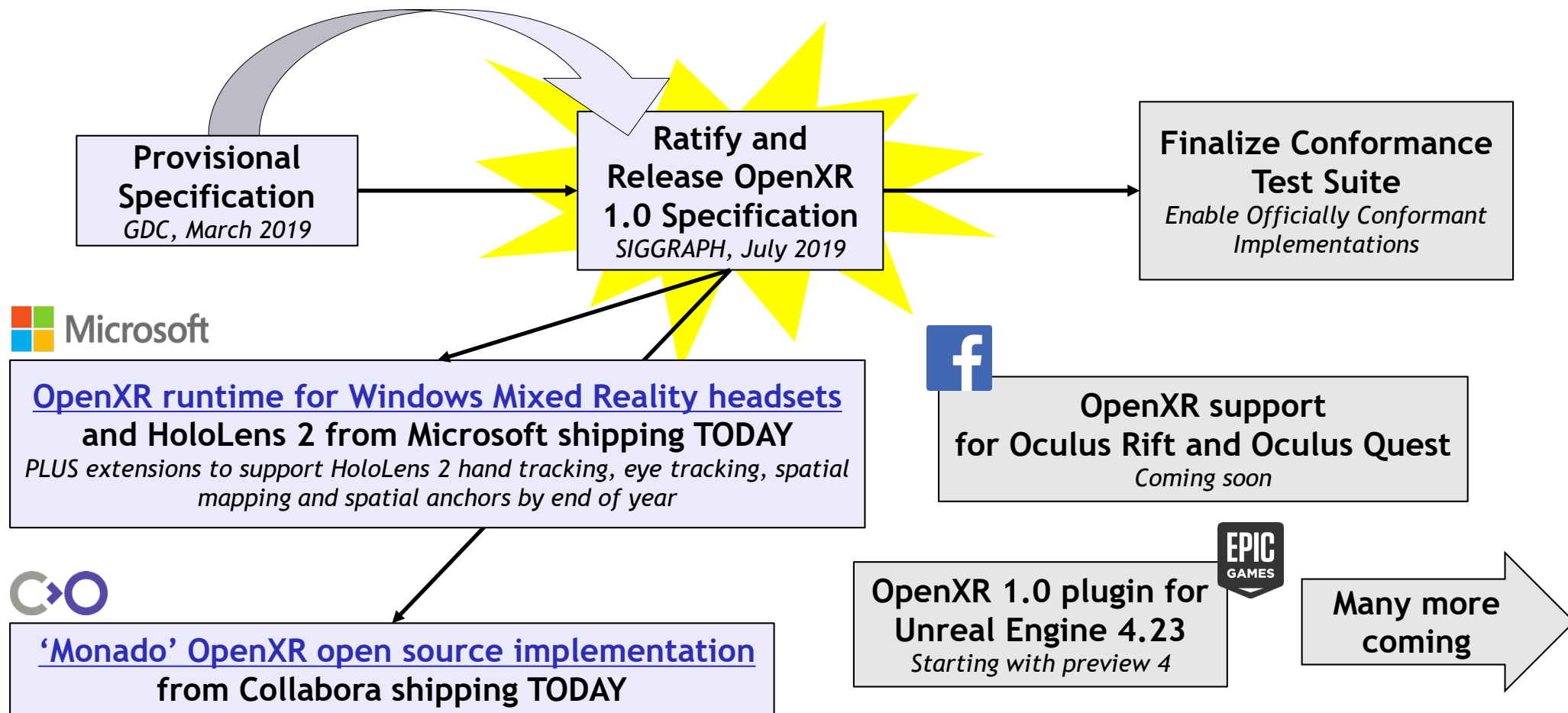
KHRONOS GROUP

# OpenXR Win-Win-Win

**XR End-Users**
Can run the apps they want on their system – reducing market confusion and increasing consumer confidence

**XR Vendors**
Can bring more applications onto their platform by leveraging the OpenXR content ecosystem

**XR ISVs**
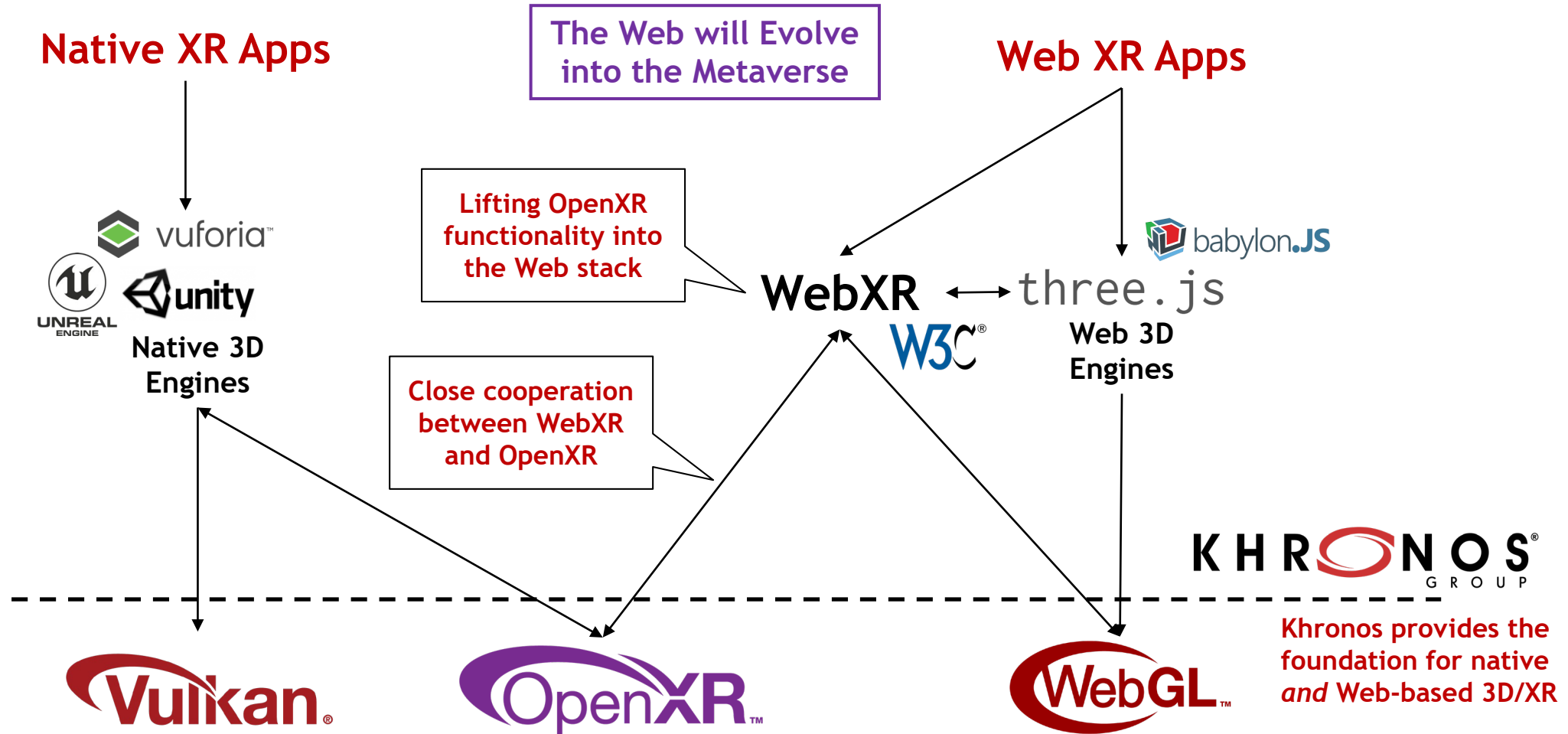Can easily ship on more platforms for increased market reach

# OpenXR 1.0 Released July 2019!

**Significant community feedback – thank you!**
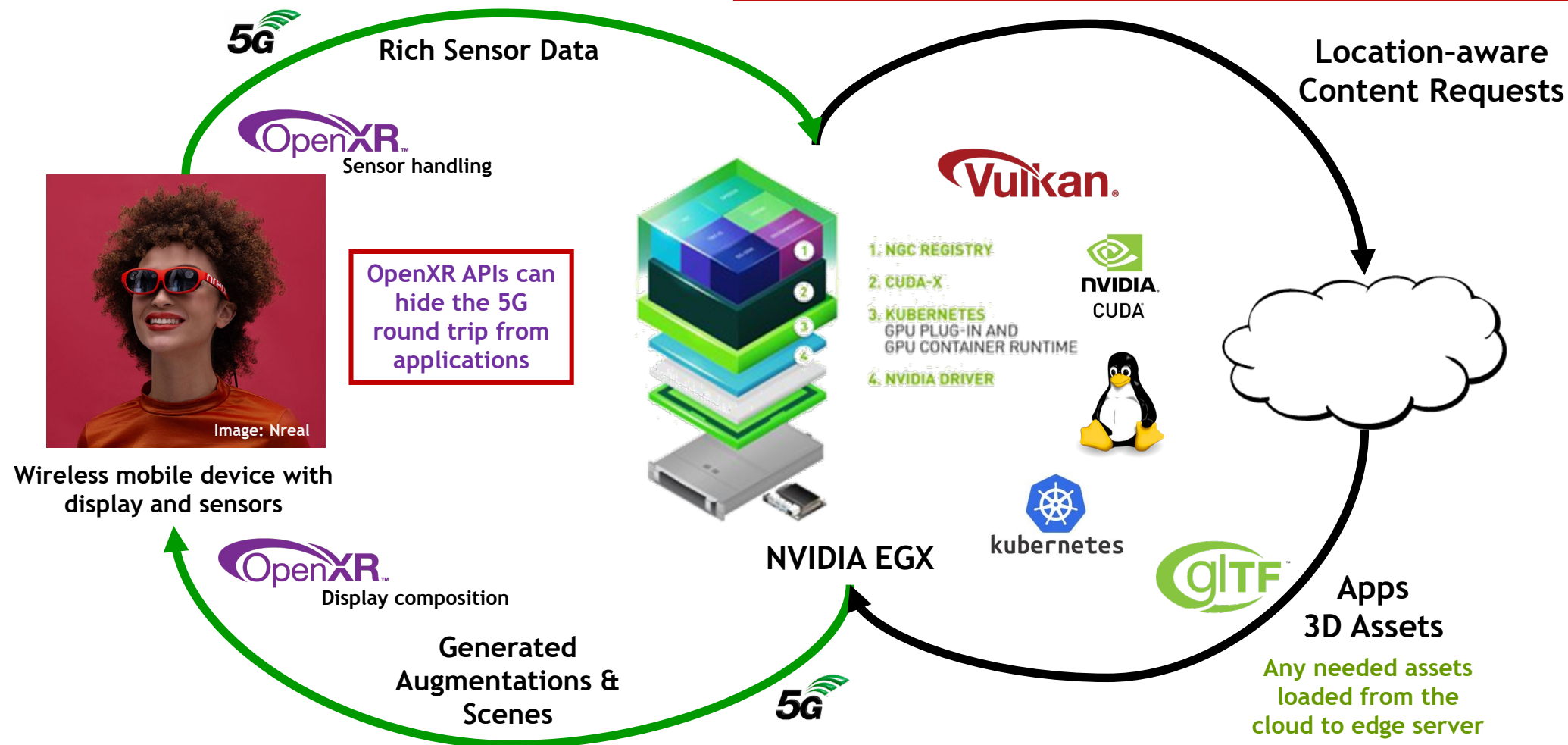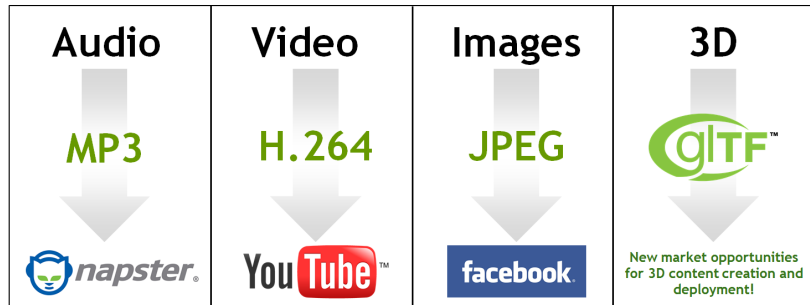Improved OpenXR input subsystem, game engine editor support, loader ...

**Provisional Specification**
*GDC, March 2019*

**Ratify and Release OpenXR 1.0 Specification**
*SIGGRAPH, July 2019*

**Finalize Conformance Test Suite**
*Enable Officially Conformant Implementations*

**Microsoft**

**OpenXR runtime for Windows Mixed Reality headsets**
**and HoloLens 2 from Microsoft shipping TODAY**
*PLUS extensions to support HoloLens 2 hand tracking, eye tracking, spatial mapping and spatial anchors by end of year*

**OpenXR support for Oculus Rift and Oculus Quest**
*Coming soon*

**'Monado' OpenXR open source implementation**
**from Collabora shipping TODAY**

**EPIC GAMES**

**OpenXR 1.0 plugin for Unreal Engine 4.23**
*Starting with preview 4*

**Many more coming**

# Bringing XR to the Web

**Native XR Apps**

The Web will Evolve into the Metaverse

**Web XR Apps**

vuforia™

Lifting OpenXR functionality into the Web stack

babylon.JS

UNREAL ENGINE  unity

**Native 3D Engines**

WebXR  ←→  three.js

W3C®

**Web 3D Engines**

Close cooperation between WebXR and OpenXR

K H R O N O S
GROUP

Vulkan®

OpenXR™

WebGL™

**Khronos provides the foundation for native _and_ Web-based 3D/XR**

# XR and 5G
## Leveraging High Bandwidth and Low Latency

MEC (Multi-access Edge Computing) Server
1. Processes sensor data, including machine learning for environmental lighting, occlusion, scene semantics, object reconstruction and UI
2. Generates imagery from 3D models, including stereo, foveal rendering, ray-tracing, optics pre-distortion, varifocal processing

**5G** Rich Sensor Data

Location-aware Content Requests

**OpenXR** Sensor handling

Image: Nreal

Wireless mobile device with display and sensors

OpenXR APIs can hide the 5G round trip from applications

**Vulkan**

1. NGC REGISTRY
2. CUDA-X
3. KUBERNETES GPU PLUG-IN AND GPU CONTAINER RUNTIME
4. NVIDIA DRIVER

NVIDIA CUDA

kubernetes

NVIDIA EGX

**OpenXR** Display composition

Generated Augmentations & Scenes

**5G**

glTF

Apps 3D Assets

Any needed assets loaded from the cloud to edge server

# glTF Real-time 3D Asset Transmission

| Audio | Video | Images | 3D |
|-------|-------|--------|-----|
| MP3 | H.264 | JPEG | glTF™ |
| napster | You Tube™ | facebook | New market opportunities for 3D content creation and deployment! |

**glTF is an efficient, reliable run-time 3D transmission format with advanced photorealistic functionality**

glTF™

Compact to Transmit ✓
Simple and Fast to Load ✓
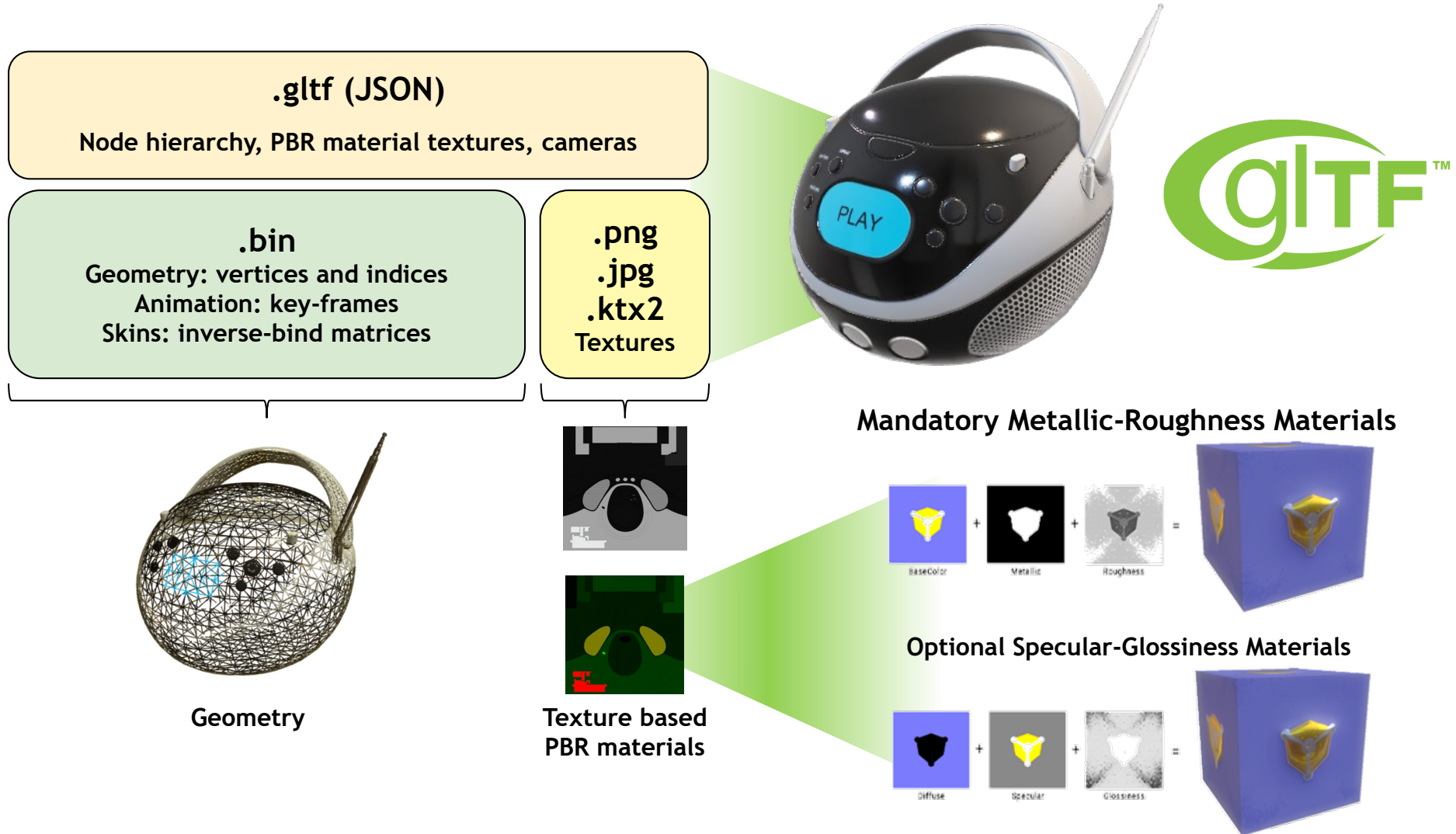Describes Full Scenes ✓
Runtime Neutral ✓
Open and Extensible ✓

WebGL™   DirectX 12   M   Vulkan.

**glTF 2.0 – June 2017
Physically Based Rendering**

## Functionality in Development

Draco Mesh Compression
Universal Compressed Textures
Second generation PBR
(absorption/attenuation, clear coat, subsurface scattering, anisotropy)
Subdivision surfaces

# glTF 2.0 Scene Description Structure



.gltf (JSON)

Node hierarchy, PBR material textures, cameras

.bin
Geometry: vertices and indices
Animation: key-frames
Skins: inverse-bind matrices

.png
.jpg
.ktx2
Textures

Geometry

Texture based
PBR materials

**Mandatory Metallic-Roughness Materials**

BaseColor + Metallic + Roughness =

**Optional Specular-Glossiness Materials**

Diffuse + Specular + Glossiness =

### Dedicated 3D Authoring Tools
MAYA · 3DS MAX · blender · Titania · Marmoset Toolbag · Paint 3D · DS SOLIDWORKS · Substance Painter · SideFX Modo · CINEMA 4D · KeyShot by Luxion

### Authoring Tools that Export 3D
COMSOL · SketchUp · Adobe Dn · MINECRAFT · Archilogic · Adobe

### VR / AR Authoring Tools
8th WALL · Microsoft Maquette BETA · spoke by moz://a make your space · UNBOUND · Oculus medium

### 3D Scanning Tools
HUAWEI 3D Live Object · scandy · eCapture 3D · Sony 3D Creator

### Convertors and Optimizers
Assimp Open Asset Import Library · CrossManager · gltfpack · OBJ2GLTF · PiXYZ SOFTWARE · FBX2glTF · SIMPLYGON · Collada2gltf · Rapid Compact · SAFE SOFTWARE

### Validation and Reference Tools
glTF Reference Viewer · gltf-vscode · AGI · glTF-asset-generator · glTF-validator · glTF-Toolkit · Microsoft

---

### Discover
**Repositories**

TURBOSQUID · Sketchfab · Poly poly.google.com

### Create
**Tools**

## glTF™ Ecosystem

### Experience
**Apps / Engines**

### Drive Demand
**Users**

Continental · web3D CONSORTIUM · NVIDIA · OGC · otoy · IKEA · EA · AMD · Bentley · Uber · shopify · TARGET · VRM CONSORTIUM

---

### Game Engines
UNREAL ENGINE · PLAYCANVAS · unity · JMonkeyEngine · GODOT Game engine · OGRE

### Web Engines
three.js · babylon.js · CLAYGL

### Apps and Engines
AUTODESK FORGE · instant3Dhub · VENTUZ · Filament · VTK · xeogl · CESIUM · 3D Builder Prep for 3D printing · STK · Qt · UX3D ENGINE · ParaView

### VR / AR Apps and Engines
magic leap · JANUSVR · hubs by moz://a · Mixed Reality Viewer · A-FRAME · worldviz · ARCore · Windows Mixed Reality Home · React 360

### Productivity and Social Apps
Office · facebook · WordPress

KHRONOS GROUP®

# glTF Ecosystem Evolution



glTF 2.0 import/export
with Blender 2.80

**Tools!**

Striving for native glTF import and Export from every tool. Catalyzed Blender IO as exemplar

glTF 2.0 – June 2017

**Consistency!**

Avoid dialects at all costs! Sample viewer and Asset Validator in open source. Sample models and asset generator for unit tests



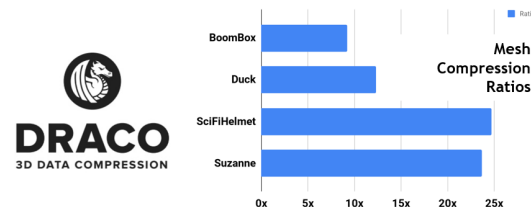Sample Viewer for accurate Ground Truth glTF renderings

**Functionality!**

Balancing functionality versus complexity. glTF is extensible – only bring widely adopted extensions into core



glTF Mesh compression extension provides up to 25x geometry compaction

glTF/Draco-enabled apps and engines

# Draco glTF Mesh Compression Extension

- **Library for compressing and decompressing 3D geometric meshes and point clouds**
  - Draco designed and built for compression efficiency and speed - great fit with glTF!
    - https://github.com/google/draco

- **Draco glTF extension launched in February 2018**
  - https://github.com/KhronosGroup/glTF/blob/master/extensions/2.0/Khronos/KHR_draco_mesh_compression/README.md

- **Google has released Draco encoders and decoders in open source**
  - C++ source code encoder to compress 3D data
  - C++ and JavaScript decoders for the encoded data
    - https://github.com/google/draco/tree/gltf_2.0_draco_extension

- **glTF/Draco compression already in use**
  - Blender, three.js, BABYLON.JS, Adobe Dimension, glTF pipeline, FBX2glTF, AMD Compressonator and glTF sample models



Mesh Compression Ratios

# Universal Textures for glTF

- **Fragmentation of GPU texture formats is significant issue for developers**
  - Binomial's 'Basis Universal' technology enables JPEG-sized texture assets
  - Transcodable on-the-fly to natively supported compressed GPU formats

- **glTF Universal Texture extension uses KTX2 subset as a flexible container**
  - Precisely defined for consistent, cross-vendor generation and validation
  - Wide range of (un)(super)compressed texture formats used in Vulkan/DirectX/Metal
  - Supports streaming and full random access to MIP levels
  - Open source tools to create, transcode and upload to WebGL, OpenGL and Vulkan
  - https://github.com/KhronosGroup/KTX-Software/tree/ktx2

Encoding decoupled from target device. One encode pass per texture asset

Transcode *on-the-fly* to a natively supported *compressed* GPU formats
Desktop: BC1-5, BC7
Mobile: ETC1/2, PVRTC1, *ASTC

Original Texture Assets (.png) → **Encode and Supercompress** → **Universal Textures** KTX2 Container with Basis Universal - compressed payload → **Transcode to GPU formats** → GPU-accelerated Texture / GPU-accelerated Texture / GPU-accelerated Texture

'toktx' OSS Tool

'libktx' OSS Tool

glTF™

*ASTC support in development

# Universal Textures: Compression Ratios



FlightHelmet_baseColor
2048 x 2048, RGB

Bytes

- **File Size**
- **GPU Size**

| Category | File Size | GPU Size |
|---|---|---|
| Uncompressed | 12,582,912 | |
| PNG | 2,778,518 | |
| JPEG | 315,619 | |
| ETC1S | 2,097,152 | |
| Basis Universal | 232,104 | |

14,000,000
10,500,000
7,000,000
3,500,000
0

JPEG must be fully decompressed into GPU memory

Universal textures can be directly transcoded to compressed GPU textures

# KTX2 and .basis files

Two complementary container formats for Basis Universal assets

**BINOMIAL**

**'Basis Universal' texture compression technology**
Enables JPG-sized textures that can be transcoded on-the-fly to natively supported *compressed* GPU formats

glTF™

**Binomial and Google open sourced 'Basis Universal' compressor and transcoder**
C++ or WebAssembly code for handling '.basis' format textures in native apps and web sites
https://github.com/binomialLLC/basis_universal
**Fine if you are in full control of your texture assets and rendering**

**Binomial's 'Basis Universal' technology contributed to glTF**
Rigorously-defined KTX2 container format supports wide range of texture formats used in Vulkan/DirectX/Metal with streaming and full random access to MIP levels
glTF extension uses KTX2 subset with supercompressed payload using Basis Universal Technology
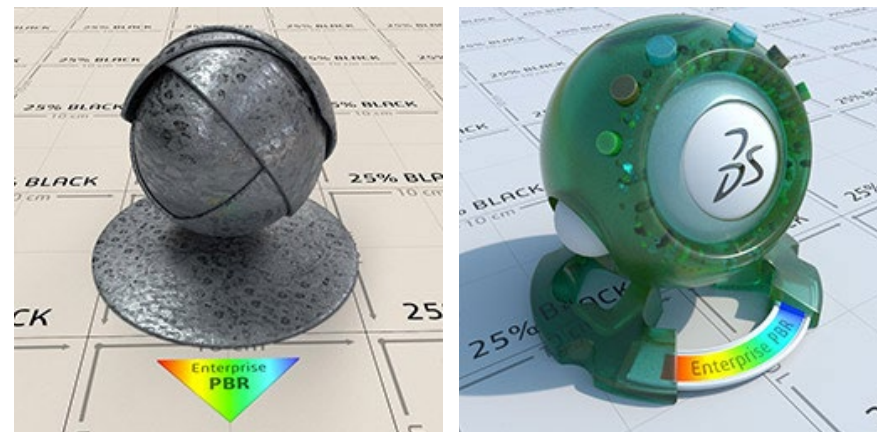**Great for cross-vendor distribution of textures to multiple applications and engines**

# Next Generation glTF PBR Materials

- **Demand for advanced PBR for photorealistic assets**
  - Beyond current 'Metallic-Roughness' and 'Specular-Glossiness'
  - E.g. Absorption/attenuation, clear coat, subsurface scattering, anisotropy

- **Extending Metallic-Roughness parameters**
  - Consistency and fallbacks for performance for any device

- **Inspiration from Dassault Systèmes Enterprise PBR Shading Model (DSPBR)**
  - https://github.com/DassaultSystemes-Technology/EnterprisePBRShadingModel/tree/master/gltf_ext

- **Wide industry collaboration for compatibility**
  - Dassault Systèmes
  - Google Filament
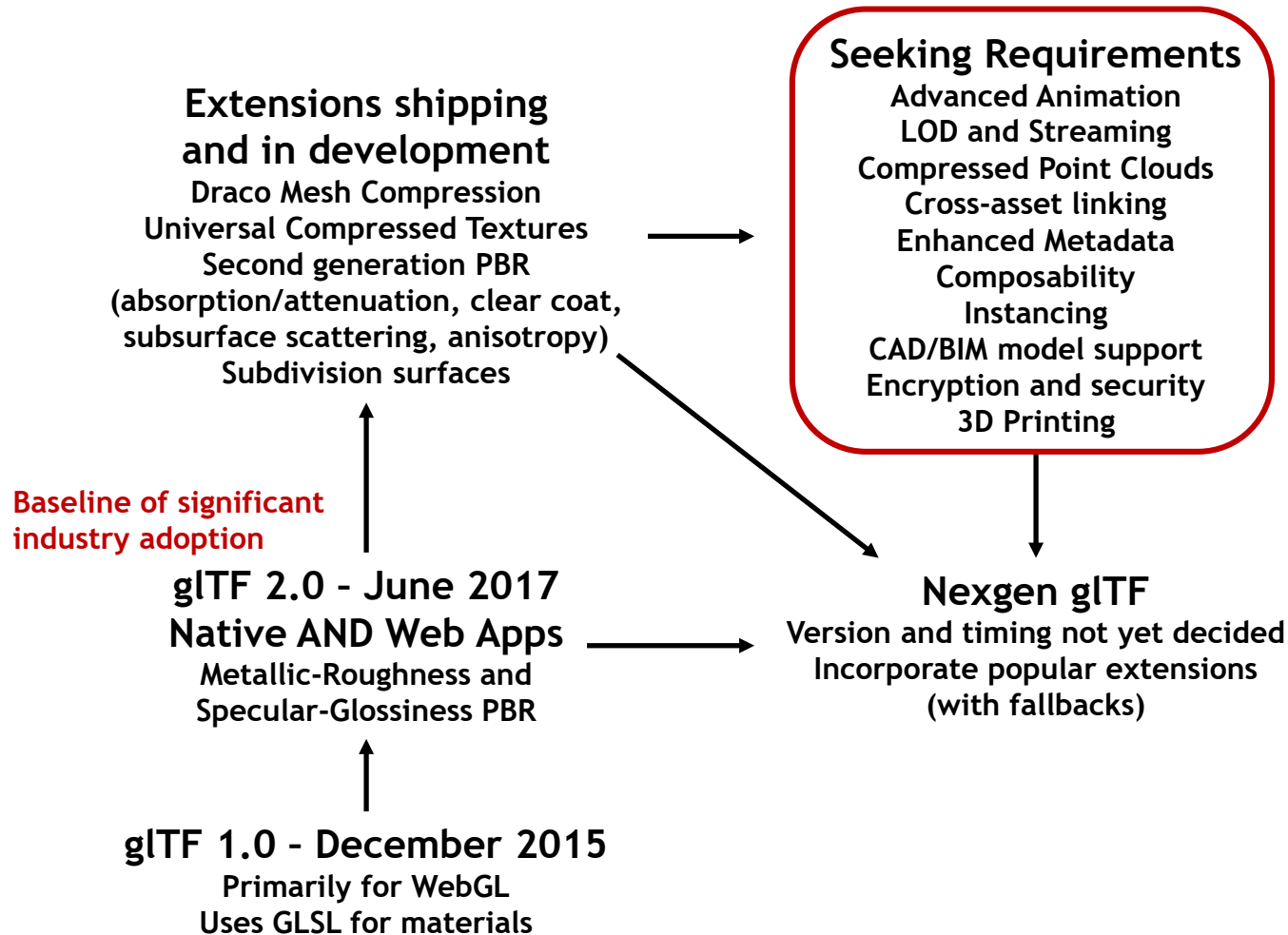  - Microsoft BabylonJS
  - NVIDIA MDL
  - OTOY Octane

Images from https://dassaultsystemes-technology.github.io/EnterprisePBRShadingModel/

**Join the GitHub Discussion!**
https://github.com/KhronosGroup/glTF/issues/1442

# glTF Evolution

**Seeking Requirements**
Advanced Animation
LOD and Streaming
Compressed Point Clouds
Cross-asset linking
Enhanced Metadata
Composability
Instancing
CAD/BIM model support
Encryption and security
3D Printing

**Extensions shipping
and in development**
Draco Mesh Compression
Universal Compressed Textures
Second generation PBR
(absorption/attenuation, clear coat,
subsurface scattering, anisotropy)
Subdivision surfaces

**The glTF Roadmap is Driven by
Developer Feedback
Join the GitHub Discussion!**
https://github.com/KhronosGroup/glTF/issues/1442

**Baseline of significant
industry adoption**

**glTF 2.0 – June 2017
Native AND Web Apps**
Metallic-Roughness and
Specular-Glossiness PBR

**Nexgen glTF**
Version and timing not yet decided
Incorporate popular extensions
(with fallbacks)

**glTF 1.0 – December 2015**
Primarily for WebGL
Uses GLSL for materials

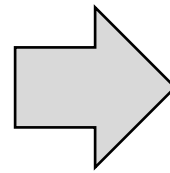# VRM Using glTF 2.0

**Hiroki Omae - Unity**

# 3D Commerce - The Opportunity

3D Commerce = E Commerce enhanced with the use of 3D Models on any platform – including VR and AR



IKEA catalog uses augmented reality to give a virtual preview of furniture in a room – August 2013

IKEA Communications AB

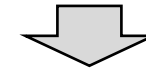**Early Experience Shows**

Increased customer engagement!
Strengthened brand loyalty!
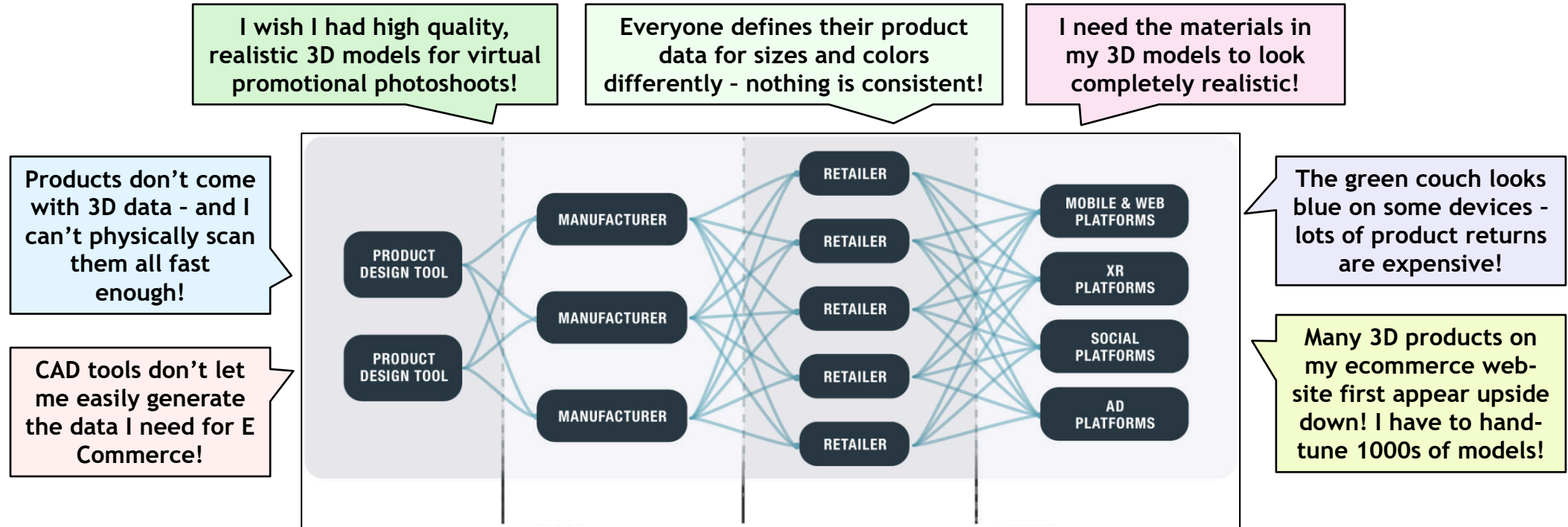Deeper product understanding!
More online sales!
Fewer returns!

**=$$$!**

**So why is 3D Commerce taking so long to become widespread?**

# 3D Commerce - Today's Reality



I wish I had high quality, realistic 3D models for virtual promotional photoshoots!

Everyone defines their product data for sizes and colors differently – nothing is consistent!

I need the materials in my 3D models to look completely realistic!

Products don't come with 3D data – and I can't physically scan them all fast enough!

CAD tools don't let me easily generate the data I need for E Commerce!

The green couch looks blue on some devices – lots of product returns are expensive!

Many 3D products on my ecommerce web-site first appear upside down! I have to hand-tune 1000s of models!

PRODUCT DESIGN TOOL · PRODUCT DESIGN TOOL · MANUFACTURER · RETAILER · MOBILE & WEB PLATFORMS · XR PLATFORMS · SOCIAL PLATFORMS · AD PLATFORMS

**Complex retail pipeline with hundreds of companies and millions of products**

**Many friction points: tooling, technical and commercial**

**3D Commerce can't reach industrial scale so...**
**Interoperability standards to the rescue!**

# Khronos 3D Commerce Initiative

**Working Group Announced SIGGRAPH 2019**



**Creating specifications and guidelines to align the 3D asset workflow from product design through manufacturing and each stage of retail to end-user delivery platforms**

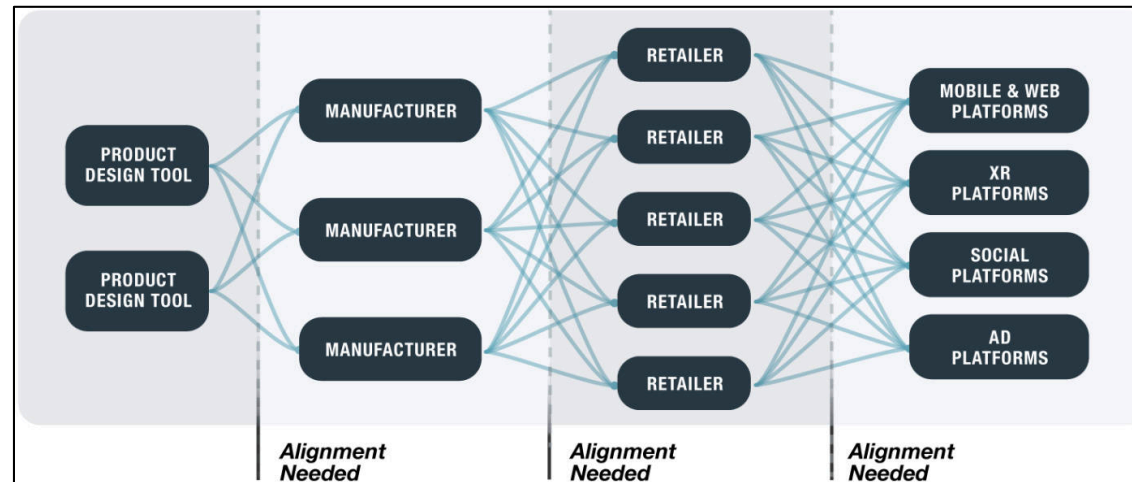**Broad Industry Participation from tooling, retail, technology and platform companies**

# 3D Commerce - Four Areas of Focus

**Asset Creation Guidelines**
For tools and product designers to create assets with consistent data to be used through the 3D Commerce pipeline

**Product Configuration**
Universal product configurability data and guidelines on how to drive consistent product display



**First Goals**
Industry cooperation to urgently develop solutions to address priority problem areas

**Metadata**
Structured metadata definitions and examples to consistently carry product information through the retail pipeline

**Viewer Validation and Certification**
Test models, reference viewer, display analysis tools and capability specifications to guarantee a consistent and accurate end user experience
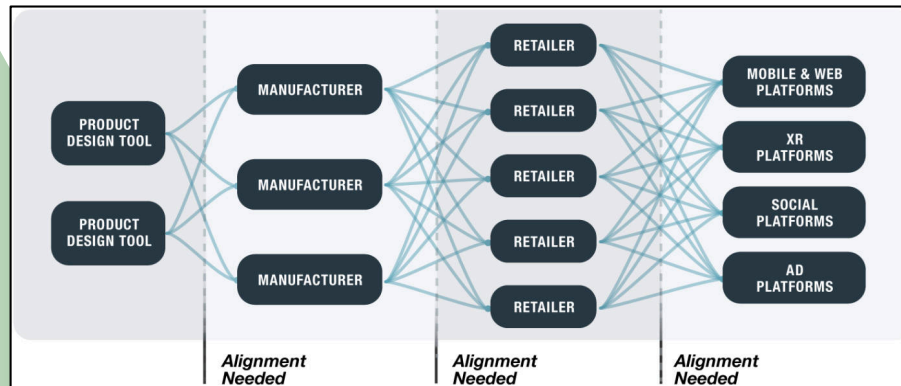
# 3D Commerce Khronos Synergy



**glTF™**
3D Asset Format

**WebGL™**
Interactive 3D on the Web

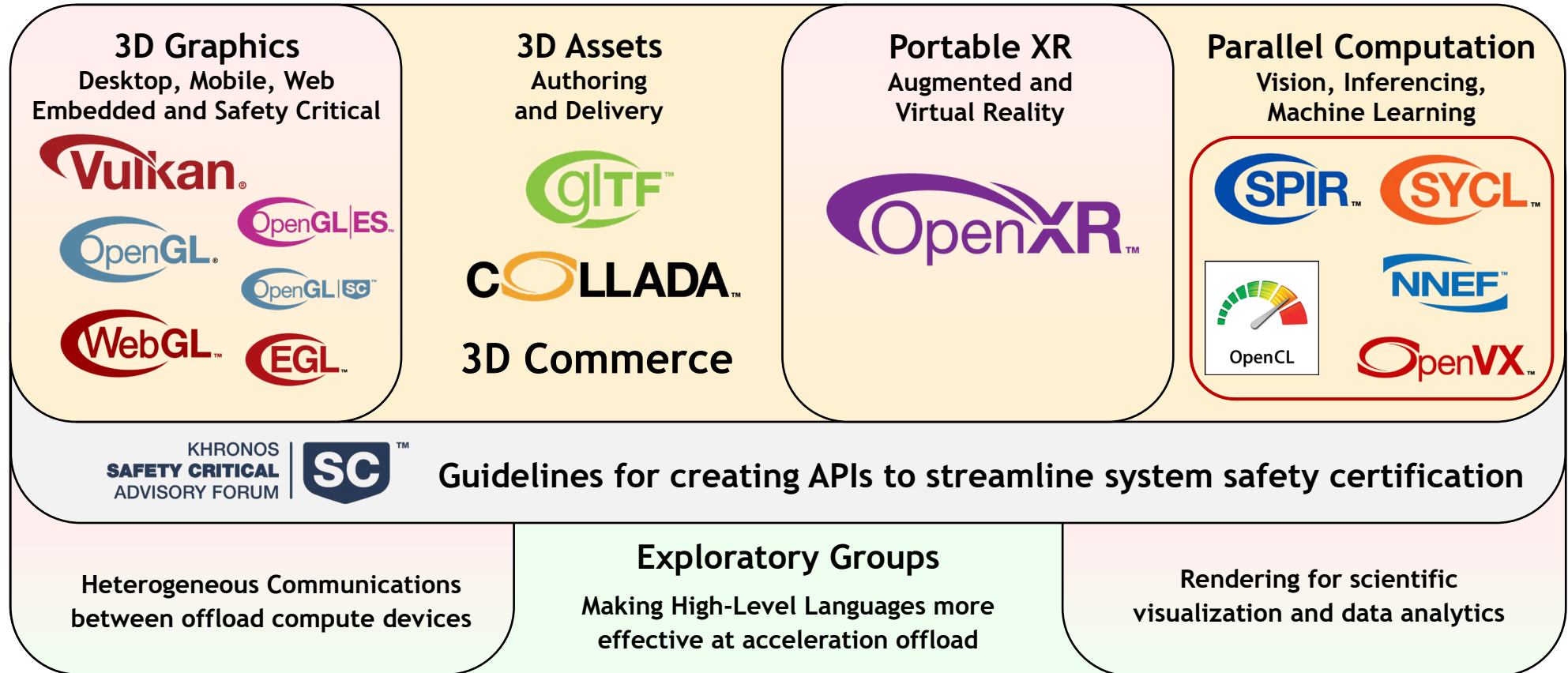**Khronos 3D Commerce**

**OpenXR™**
Portable AR and VR Apps

**Vulkan®**
High-performance
cross-platform 3D graphics

**OpenVX™**
Vision processing and
inferencing for AR and scanning

# Khronos Active Initiatives

**3D Graphics**
Desktop, Mobile, Web
Embedded and Safety Critical

**Vulkan.**
**OpenGL|ES.**
**OpenGL.**
**OpenGL|SC.**
**WebGL.**
**EGL.**

**3D Assets**
Authoring
and Delivery

**glTF.**
**COLLADA.**

**3D Commerce**

**Portable XR**
Augmented and
Virtual Reality

**OpenXR.**

**Parallel Computation**
Vision, Inferencing,
Machine Learning

**SPIR.**
**SYCL.**
OpenCL
**NNEF.**
**OpenVX.**

KHRONOS
**SAFETY CRITICAL**
ADVISORY FORUM **SC** ™
Guidelines for creating APIs to streamline system safety certification

Heterogeneous Communications
between offload compute devices

**Exploratory Groups**
Making High-Level Languages more
effective at acceleration offload

Rendering for scientific
visualization and data analytics

# Khronos Open Standard Compute APIs

**High-level APIs**

**SYCL**™

Single source C++ Programming
with Compute Acceleration

**OpenVX**™

Vision and Inferencing
Acceleration

Import ←

**NNEF**™

Trained Neural Network
Exchange Format

**Low-level APIs**

**Vulkan**®

GPU Rendering +
Compute Acceleration

↓

| GPU |

OpenCL

Heterogeneous Compute
Acceleration

↓

| CPU | GPU |
| FPGA | DSP |
| Custom Hardware | |

**Increasing interest in parallel heterogonous compute acceleration to combat the 'End of Moore's Law'**

**KHRONOS** GROUP

# SYCL Single Source C++ Parallel Programming

**SYCL is ideal for accelerating larger C++-based engines and applications**

**Multiple SYCL libraries for vision and inferencing SYCL-BLAS, SYCL-DNN, SYCL-Eigen, SYCL Parallel STL**

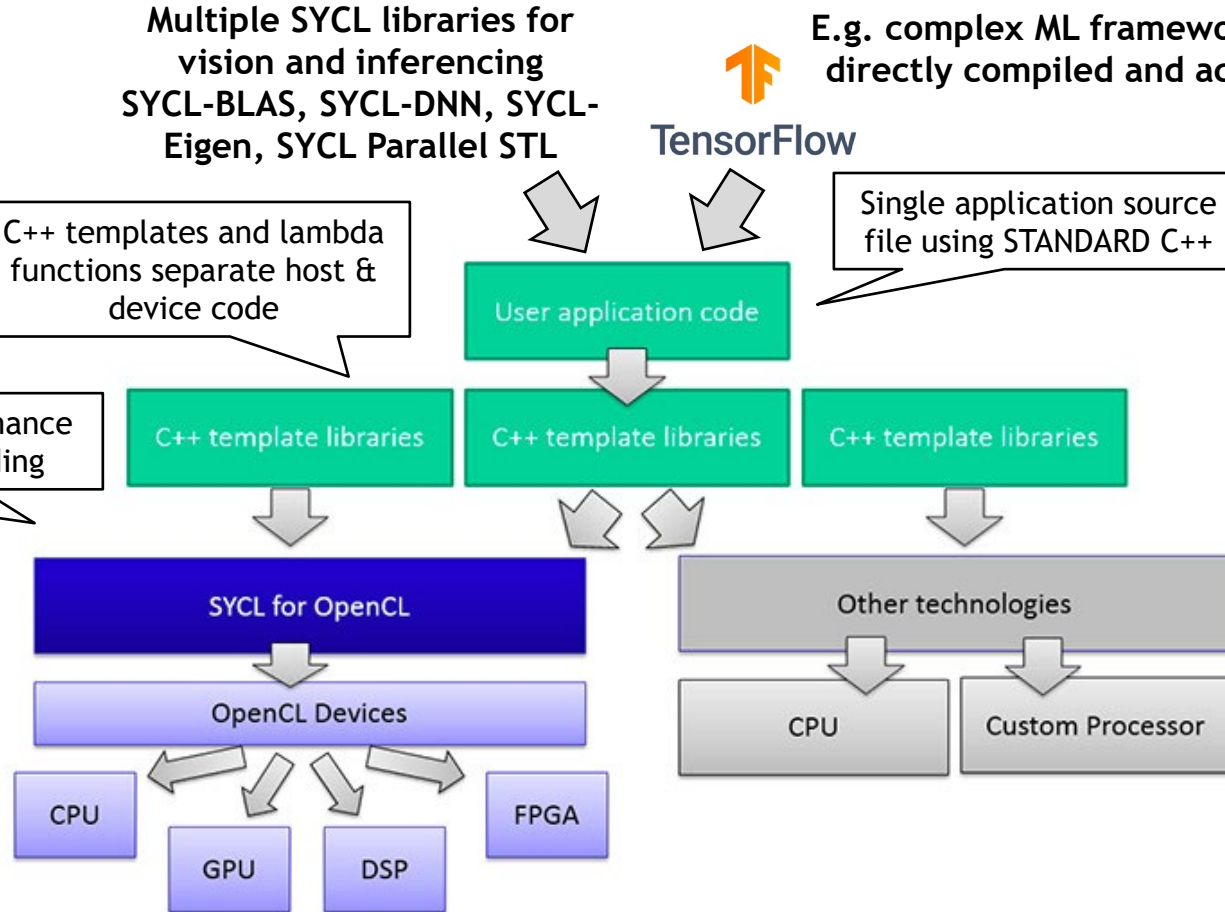**E.g. complex ML frameworks can be directly compiled and accelerated**

TensorFlow

C++ templates and lambda functions separate host & device code

Single application source file using STANDARD C++

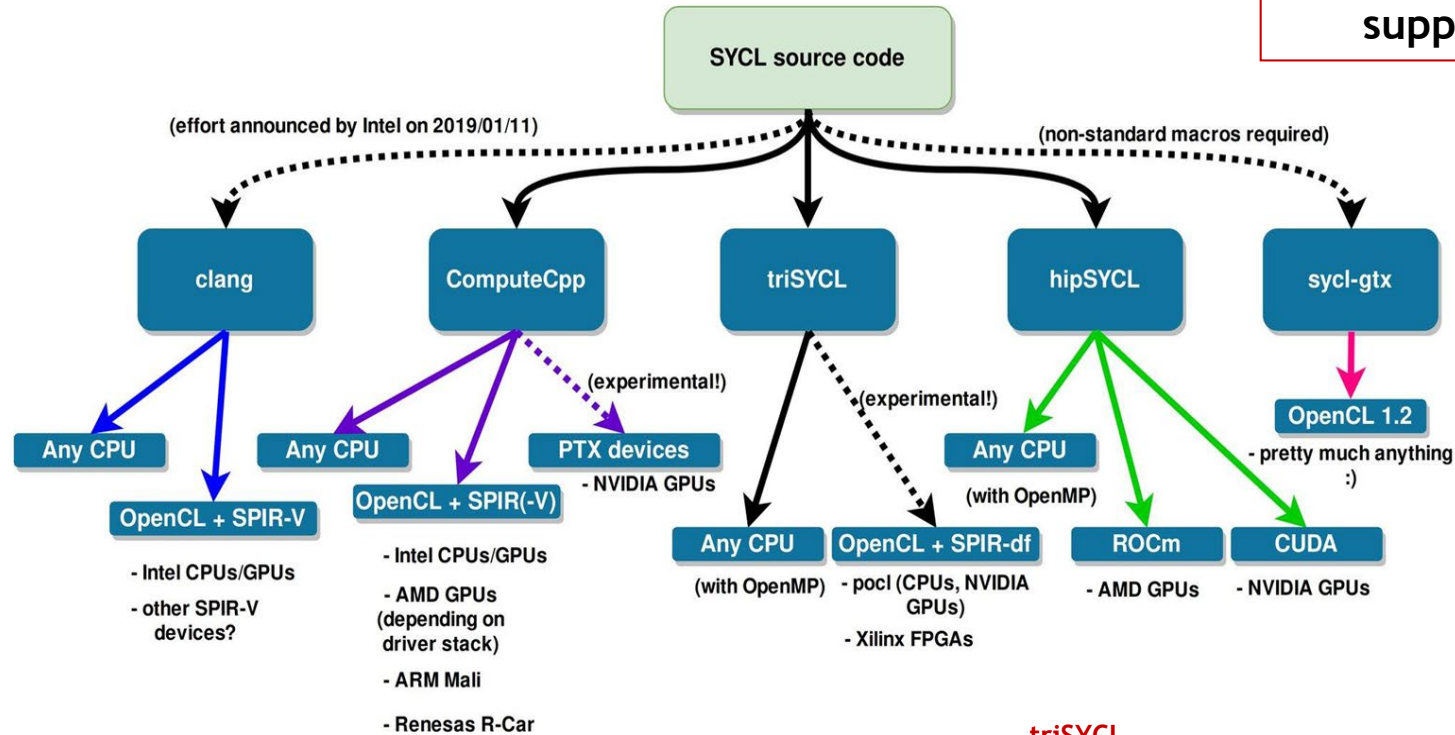C++ Kernel Fusion can give better performance on complex apps and libs than hand-coding

User application code

C++ template libraries

C++ template libraries

C++ template libraries

SYCL

OpenCL

SYCL for OpenCL

Other technologies

OpenCL Devices

CPU

Custom Processor

Accelerated code passed into device OpenCL compilers

CPU

GPU

DSP

FPGA

# SYCL Implementations

SYCL enables Khronos to influence ISO to enable standard C++ to (eventually) support heterogenous compute

**SYCL source code**

(effort announced by Intel on 2019/01/11)

(non-standard macros required)

**clang**
**ComputeCpp**
**triSYCL**
**hipSYCL**
**sycl-gtx**

(experimental!)

**Any CPU**

**Any CPU**

**PTX devices**
- NVIDIA GPUs

(experimental!)

**Any CPU**
(with OpenMP)

**OpenCL 1.2**
- pretty much anything :)

**OpenCL + SPIR-V**
- Intel CPUs/GPUs
- other SPIR-V devices?

**OpenCL + SPIR(-V)**
- Intel CPUs/GPUs
- AMD GPUs (depending on driver stack)
- ARM Mali
- Renesas R-Car

**Any CPU**
(with OpenMP)

**OpenCL + SPIR-df**
- pocl (CPUs, NVIDIA GPUs)
- Xilinx FPGAs

**ROCm**
- AMD GPUs

**CUDA**
- NVIDIA GPUs

## Multiple Backend Support Coming
SYCL beginning to be supported on low-level APIs in addition to OpenCL e.g. Vulkan and CUDA
http://sycl.tech

## Intel Adoption
Intel's 'One API' Initiative uses SYCL
https://newsroom.intel.com/news/intels-one-api-project-delivers-unified-programming-model-across-diverse-architectures/#gs.bydj6z

**LLVM/clang SYCL Compiler**
Compiles C++-based SYCL source files into code for both CPU and a wide range of compute accelerators

**ComputeCpp**
Codeplay Software's v1.2.1 conformant implementation available to download today

**triSYCL**
Open-source test-bed to experiment with the specification of the OpenCL SYCL C++ layer and to give feedback to Khronos

**HipSYCL**
SYCL 1.2.1 implementation that builds upon NVIDIA CUDA/AMD HIP/ROCm
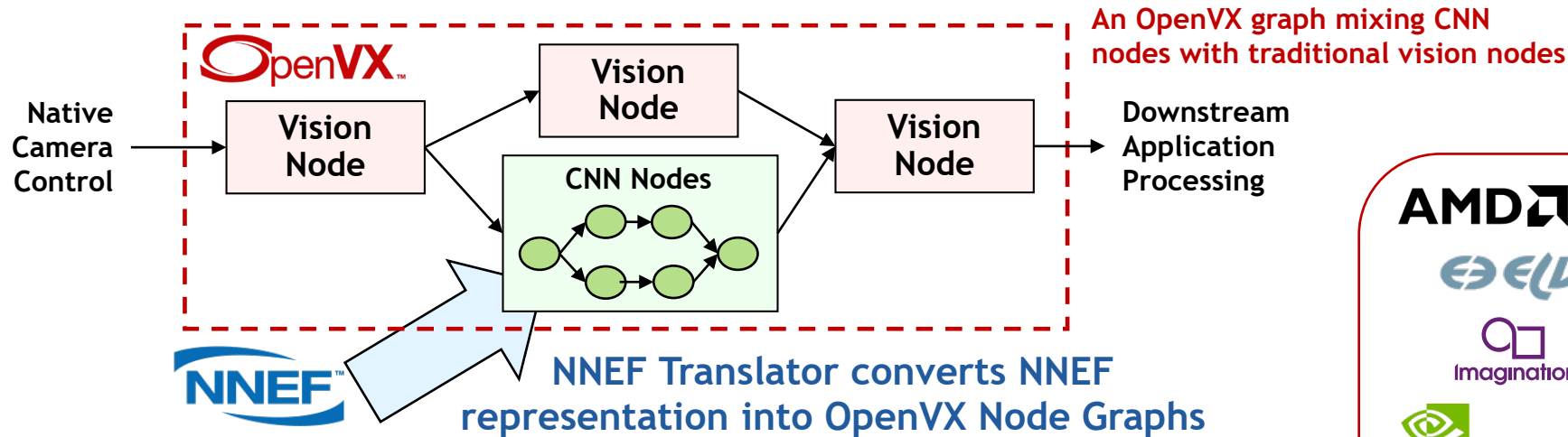
# OpenVX Cross-Vendor Inferencing

## OpenVX

A high-level graph-based abstraction for portable, efficient vision processing
Optimized OpenVX drivers created and shipped by processor vendors
Can be implemented on almost any hardware or processor
Graph can contain vision processing and NN nodes – enables global optimizations
Run-time graph execution can be almost completely autonomously from the host CPU



An OpenVX graph mixing CNN nodes with traditional vision nodes

NNEF Translator converts NNEF representation into OpenVX Node Graphs

Hardware Implementations

## Performance comparable to hand-optimized, non-portable code

Real, complex applications on real, complex hardware
Much lower development effort than hand-optimized code

# OpenVX 1.3 Released October 2019


### OpenVX
Version 1.3

**OpenVX 1.3 Feature Sets**

Enables deployment flexibility while avoiding fragmentation
Implementations with one or more complete feature sets are conformant
- Baseline Graph Infrastructure (enables other Feature Sets)
- Default Vision Functions
- Enhanced Vision Functions (introduced in OpenVX 1.2)
- Neural Network Inferencing (including tensor objects)
- NNEF Kernel import (including tensor objects)
- Binary Images
- Safety Critical (reduced features for easier safety certification)
https://www.khronos.org/registry/OpenVX/specs/1.3/html/OpenVX_Specification_1_3.html

**Open Source Prototype OpenVX 1.3
Conformance Test Suite**
Finalization expected before the end of 2019
https://github.com/KhronosGroup/OpenVX-cts/tree/openvx_1.3

**Open Source OpenVX Tutorial
and Code Samples**
https://github.com/rgiduthuri/openvx_tutorial

**Open source OpenVX 1.3 for Raspberry Pi**

Raspberry Pi 3 Model B with Raspbian OS
Automatic optimization of memory access patterns via tiling and chaining
Highly optimized kernels leveraging multimedia instruction set
Automatic parallelization for multicore CPUs and GPUs
Automatic merging of common kernel sequences
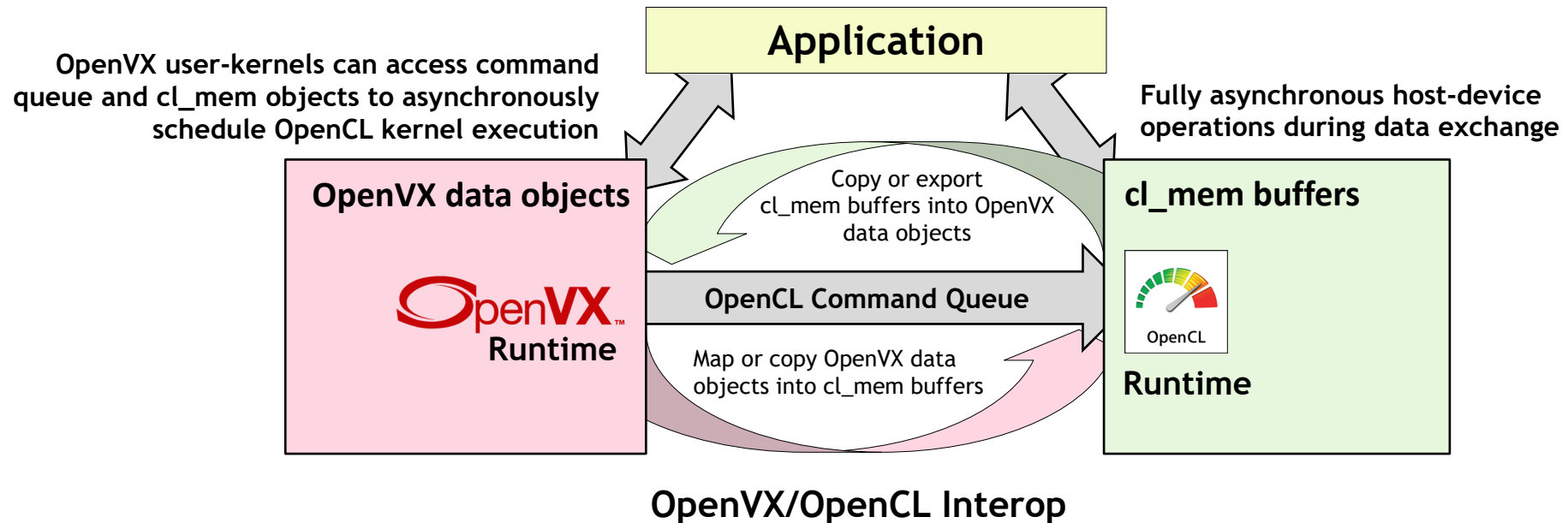https://github.com/KhronosGroup/OpenVX-sample-impl/tree/openvx_1.3

# Extending OpenVX with Custom Nodes

## OpenVX/OpenCL Interop
- Provisional Extension
- Enables custom OpenCL acceleration to be invoked from OpenVX User Kernels
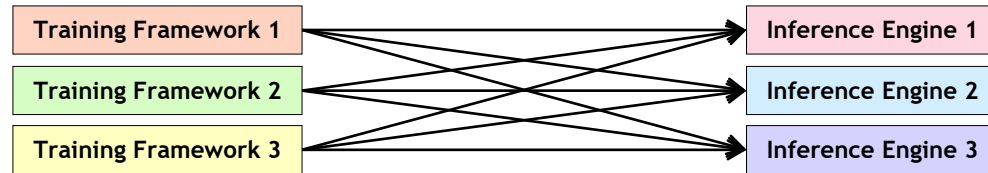- Memory objects can be mapped or copied

## Kernel/Graph Import
- Provisional Extension
- Defines container for executable or IR code
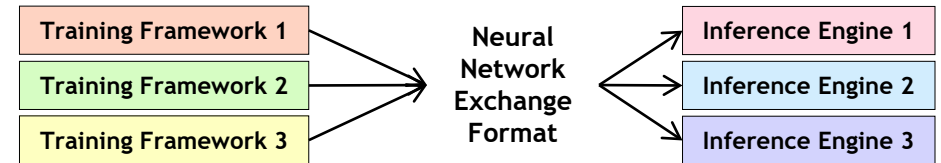- Enables arbitrary code to be inserted as an OpenVX Node in a graph

**Application**

OpenVX user-kernels can access command queue and cl_mem objects to asynchronously schedule OpenCL kernel execution

Fully asynchronous host-device operations during data exchange

**OpenVX data objects**

OpenVX™ Runtime

**cl_mem buffers**

OpenCL

Runtime

Copy or export cl_mem buffers into OpenVX data objects

**OpenCL Command Queue**

Map or copy OpenVX data objects into cl_mem buffers

**OpenVX/OpenCL Interop**

# Neural Network Exchange Formats

## Before - Training and Inferencing Fragmentation

| Training Framework 1 | Training Framework 2 | Training Framework 3 |

→ Inference Engine 1, Inference Engine 2, Inference Engine 3

**Every Inferencing Engine needs a custom importer from every Framework**

## After - NN Training and Inferencing Interoperability

Training Framework 1, Training Framework 2, Training Framework 3 → **Neural Network Exchange Format** → Inference Engine 1, Inference Engine 2, Inference Engine 3

**Common Optimization and processing tools**

## Two Neural Network Exchange Format Initiatives

| NNEF | ONNX |
|---|---|
| Embedded Inferencing Import | Training Interchange |
| Defined Specification | Open Source Project |
| Multi-company Governance at Khronos | Initiated by Facebook & Microsoft |
| Stability for hardware deployment | Software stack flexibility |

**ONNX and NNEF are Complementary**
ONNX moves quickly to track authoring framework updates
NNEF provides a stable bridge from training into edge inferencing engines

# NNEF and ONNX Industry Support

## NNEF V1.0 released in August 2018
After positive industry feedback on Provisional Specification.
Maintenance update issued in September 2019
Extensions to V1.0 released for expanded functionality



**NNEF Working Group Participants**
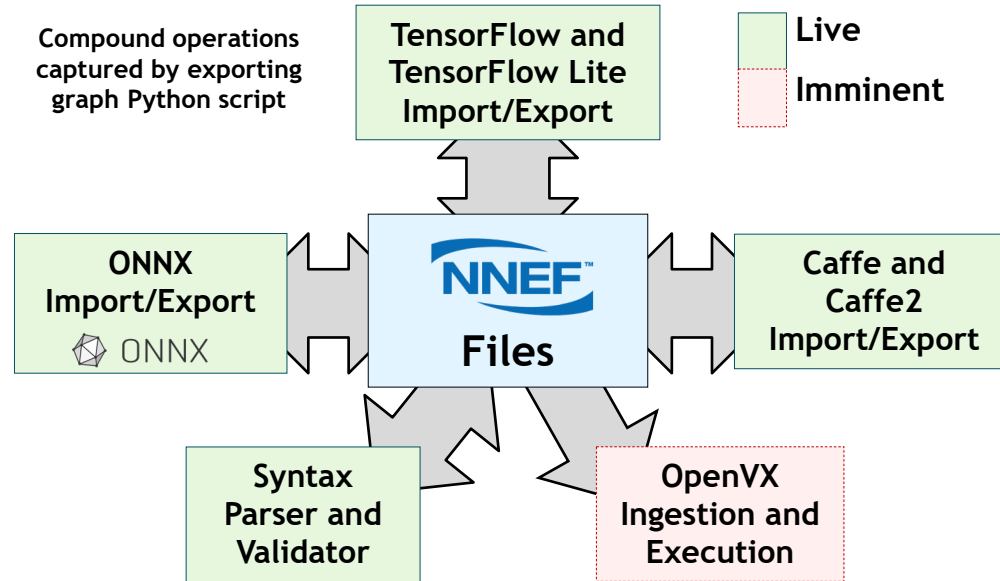
## ONNX 1.6 Released in September 2019
Introduced support for Quantization
ONNX Runtime being integrated with GPU inferencing
engines such as NVIDIA TensorRT



**ONNX Supporters**

# NNEF Tools Ecosystem

Compound operations
captured by exporting
graph Python script

TensorFlow and
TensorFlow Lite
Import/Export

Live

Imminent

ONNX
Import/Export
ONNX

**NNEF™**
Files

Caffe and
Caffe2
Import/Export

Syntax
Parser and
Validator

OpenVX
Ingestion and
Execution

NNEF open source projects hosted on Khronos
NNEF GitHub repository under Apache 2.0
https://github.com/KhronosGroup/NNEF-Tools

**NNEF™**

## NNEF Model Zoo
Now available on GitHub. Useful for
checking that ingested NNEF produces
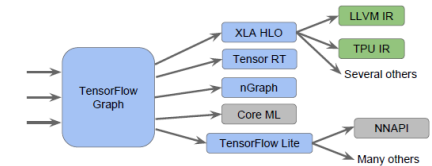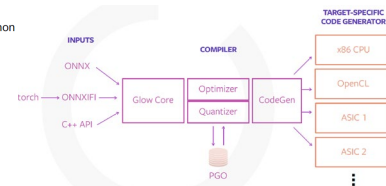acceptable results on target system

## NNEF adopts a rigorous approach to design lifecycle
Especially important for safety-critical or
mission-critical applications in automotive,
industrial and infrastructure markets

# Primary Machine Learning Compilers

| | tvm.ai (amazon) | plaidML (intel) | Glow (Facebook) | XLA (Google) |
|---|---|---|---|---|
| **Import Formats** | Caffe, Keras, MXNet, ONNX | TensorFlow Graph, MXNet, PaddlePaddle, Keras, ONNX | PyTorch, ONNX | TensorFlow Graph, PyTorch |
| **Front-end / IR** | NNVM / Relay IR | nGraph / Stripe IR | Glow Core / Glow IR | XLA HLO  MLIR |
| **Output** | OpenCL, LLVM, CUDA, Metal | OpenCL, LLVM, CUDA | OpenCL LLVM | LLVM, TPU IR, XLA IR TensorFlow Lite / NNAPI (inc. HW accel) |

# ML Compiler Steps

**Embedded NN Compilers**
CEVA Deep Neural Network (CDNN)
Cadence Xtensa Neural Network Compiler (XNNC)

## Consistent Steps

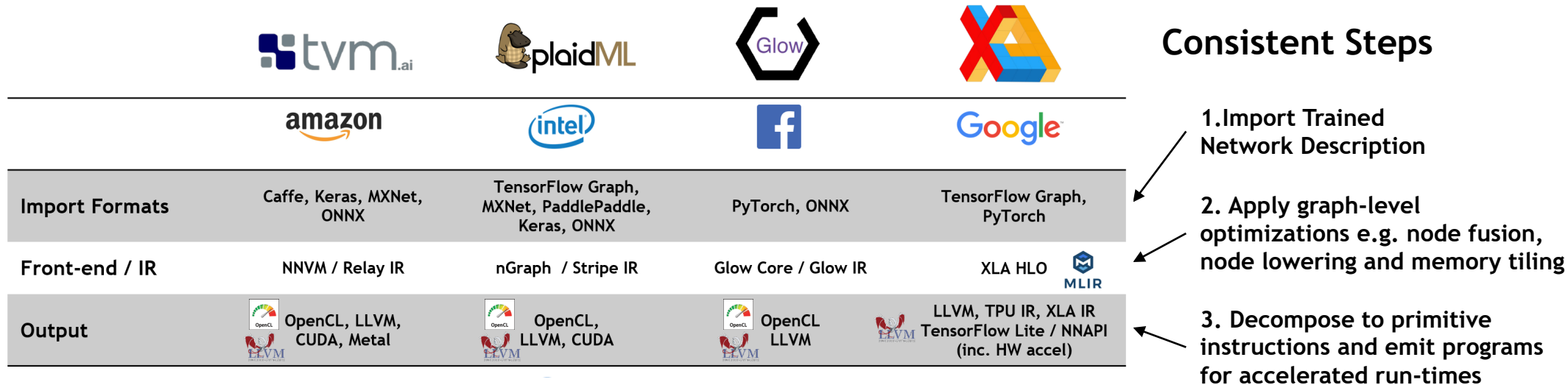|  | tvm.ai | plaidML | Glow | XLA |
|---|---|---|---|---|
|  | amazon | (intel) | facebook | Google |
| **Import Formats** | Caffe, Keras, MXNet, ONNX | TensorFlow Graph, MXNet, PaddlePaddle, Keras, ONNX | PyTorch, ONNX | TensorFlow Graph, PyTorch |
| **Front-end / IR** | NNVM / Relay IR | nGraph / Stripe IR | Glow Core / Glow IR | XLA HLO  MLIR |
| **Output** | OpenCL, LLVM, CUDA, Metal | OpenCL, LLVM, CUDA | OpenCL LLVM | LLVM, TPU IR, XLA IR TensorFlow Lite / NNAPI (inc. HW accel) |

1. Import Trained Network Description

2. Apply graph-level optimizations e.g. node fusion, node lowering and memory tiling

3. Decompose to primitive instructions and emit programs for accelerated run-times
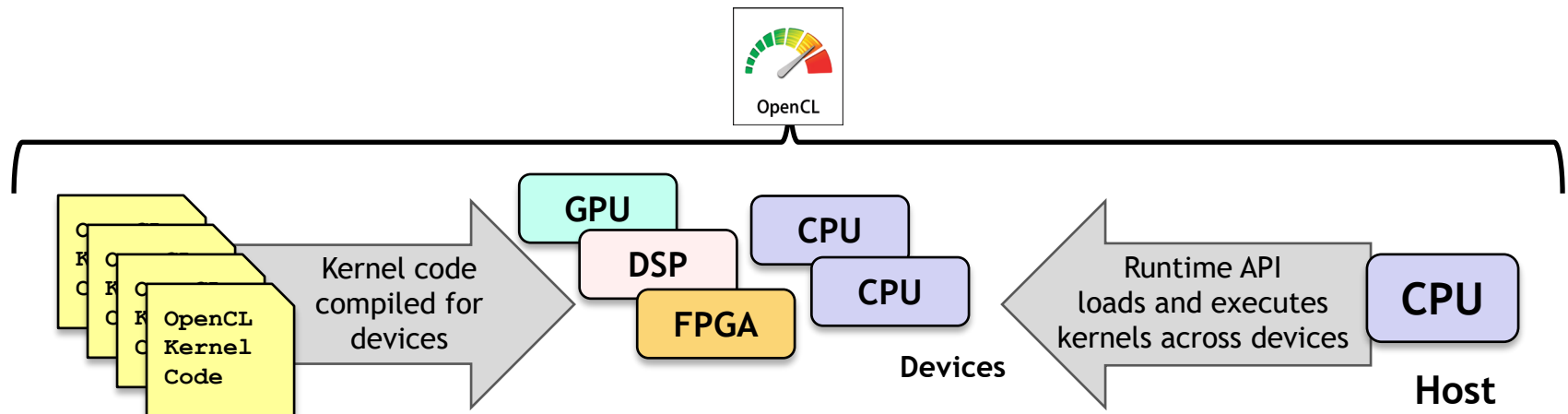
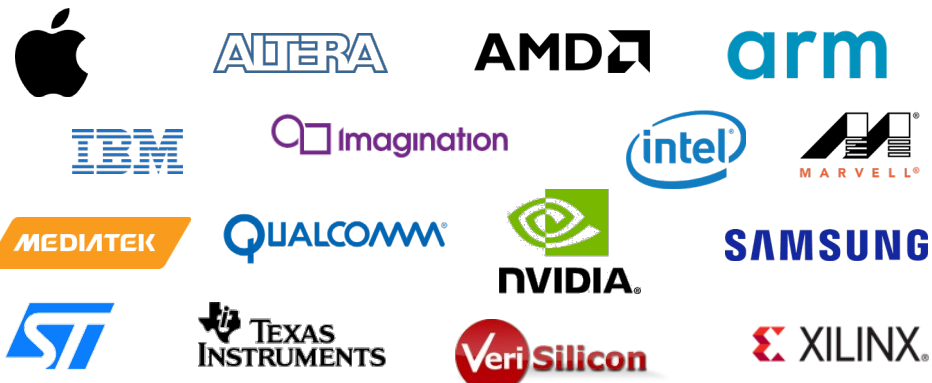**Fast progress but still area of intense research**
If compiler optimizations are effective - hardware accelerator APIs can stay 'simple' and won't need complex metacommands (combined primitive commands) like DirectML

# OpenCL - Low-level Parallel Programing

- **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

| OpenCL Kernel Code | Kernel code compiled for devices | GPU / DSP / FPGA / CPU / CPU | Runtime API loads and executes kernels across devices | CPU |

Devices

Host

# OpenCL is Widely Deployed and Used

**Hardware Implementations**

Apple, ALTERA, AMD, arm, IBM, Imagination, intel, MARVELL, MEDIATEK, QUALCOMM, NVIDIA, SAMSUNG, ST, TEXAS INSTRUMENTS, VeriSilicon, XILINX

**Desktop Creative Apps**

F, Adobe, blender, Capture One, otoy, DASSAULT SYSTEMES, SONY, Modo, ArcSoft, SideFX, CyberLink, CHAOSGROUP, ptc, Blackmagicdesign, GIMP, AUTODESK

**Linear Algebra Libraries**

CLBlast, ViennaCL, SYCL-BLAS

**Parallel Computation Languages**

OpenACC DIRECTIVES FOR ACCELERATORS, SYCL, aparapi, PyOpenCL

**Math and Physics Libraries**

ArrayFire C, C++, Fortran, MATHLAB, Wolfram Mathematica, GNU Octave, BULLET PHYSICS LIBRARY

**Vision and Imaging Libraries**

Image Magick, VisionCpp, Halide, OpenVX, OpenCV

**Machine Learning Inferencing Compilers**

Glow, tvm, plaidML

**Machine Learning Libraries**

OPENVINO Intel, Huawei MACE Mobile AI Compute Engine, Intel clDNN, Synopsis MetaWare EV, Android NNAPI, Qualcomm Neural Processing SDK for AI, Caffe, SYCL-DNN, Arm Compute Library, TI Deep Learning Library (TIDL), Acuity ML/AI Software Foundation VeriSilicon

# OpenCL Evolution



**OpenCL Extension Specs**
Scratch-Pad Memory Management
Vulkan / OpenCL Interop
Extended Subgroups
SPIR-V 1.4 ingestion for compiler efficiency
SPIR-V Extended debug info

**Integration of Extensions
plus New Core functionality**

OpenCL

May 2017
OpenCL 2.2

OpenCL

Target 2020
'OpenCL Next'

**Focus for OpenCL Next is
'Deployment Flexibility'**
Flexible Profile enables embedded vendors
to ship targeted functionality for their
customers and be officially conformant

**Regular Maintenance Updates**
Regular updates for spec clarifications,
formatting and bug fixes
https://www.khronos.org/registry/OpenCL/

**Repeat Cycle for next
Core Specification**

# Deploying OpenCL C Over Vulkan

- **Clspv – Google's experimental compiler for OpenCL C to Vulkan SPIR-V**
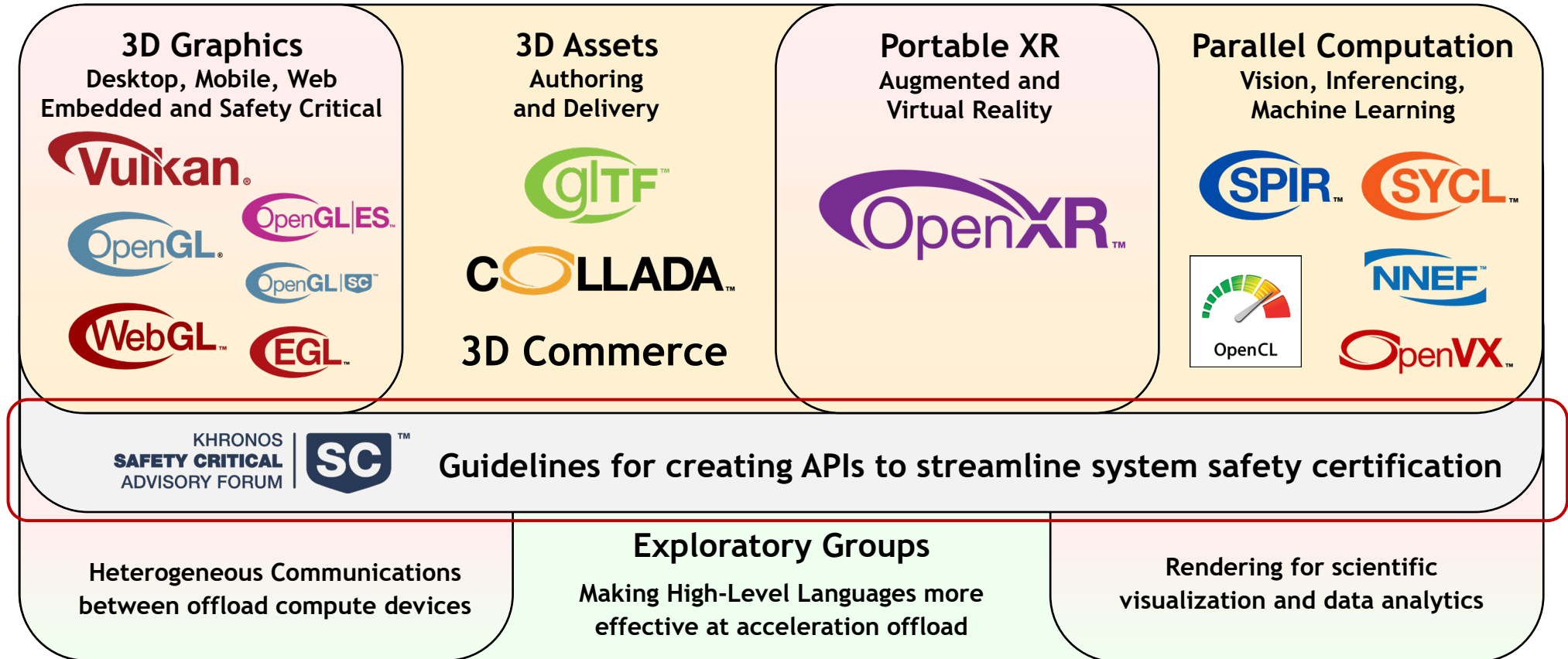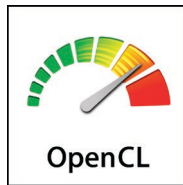  - Open source - tracks top-of-tree LLVM and clang, not a fork

- **Adobe Premiere Rush has 200K lines of OpenCL C kernel code**
  - Professional-quality, cross-platform video capture and editing system
  - Now shipping on Android on Vulkan



Adobe Premiere Rush — Video Editor
Adobe   Video Players & Editors           ★★★★★ 871
Everyone
Offers in-app purchases
This app is compatible with your device.
Add to Wishlist        Install

Shoot, edit, and share online videos anywhere.

**OpenCL C Source**

**OpenCL Host Code**

Google

**Clspv Compiler**

**Run-time API Translator**

**Prototype open source project**
https://github.com/google/clspv

SPIR™

**Runtime**

Vulkan®

**Prototype open source project**
https://github.com/kpet/clvk

# Khronos Active Initiatives

# Safety Critical GPU API Evolution



**OpenGL SC 1.0 - 2005**
Fixed function graphics safety critical subset

**OpenGL SC 2.0 - April 2016**
Programmable Shaders safety critical subset

**New Generation Safety Critical APIs for Graphics, Compute and Display**

**OpenGL ES 1.0 - 2003**
Fixed function graphics

**OpenGL ES 2.0 - 2007**
Programmable Shaders

**Vulkan 1.0 - 2016**
Explicit Graphics and Compute

**OpenCL**

**Potential OpenCL SC work will leverage the deployment flexibility of 'OpenCL Next' to minimize API surface area**

Rendering    Compute    Display

**Industry Need**
for GPU Acceleration APIs designed to ease system safety certification is increasing
ISO 26262 / ASIL-D

**ISO 26262**

**Vulkan is Compelling Starting Point for SC GPU API Design**
- Widely adopted, royalty-free open standard
- Low-level explicit API - smaller surface area than OpenGL
  - Not burdened by debug functionality
  - Very little internal state
  - Well-defined thread behavior

**Clearly Definable Design Goals to Adapt Vulkan for SC**
Reduce driver size and complexity
-> Offline pipeline creation, no dynamic display resolutions
Deterministic Behavior
-> No ignored parameters, static memory management, eliminate undefined behaviors
Robust Error Handling
-> Error callbacks so app can respond, Fatal error callbacks for fast recovery initiation
C API - MISRA C Compliance

**Vulkan|SC**

**Khronos Vulkan SC Working Group started work in February 2019**

# Khronos Standards Immersive Computing



**Download 3D object and scene data**

**Vision and sensor processing - including neural network inferencing for machine learning**

**High-performance, low-latency 3D Graphics**

**Portable interaction with VR/AR sensor, haptic and display devices**

# Khronos Proven Process and Organization

Open membership.
Any company is welcome to join.
One company one vote

Open specifications.
ROYALTY-FREE through a strong,
modern IP Framework

Any member, or non-member, can propose new standards initiatives

**Software**

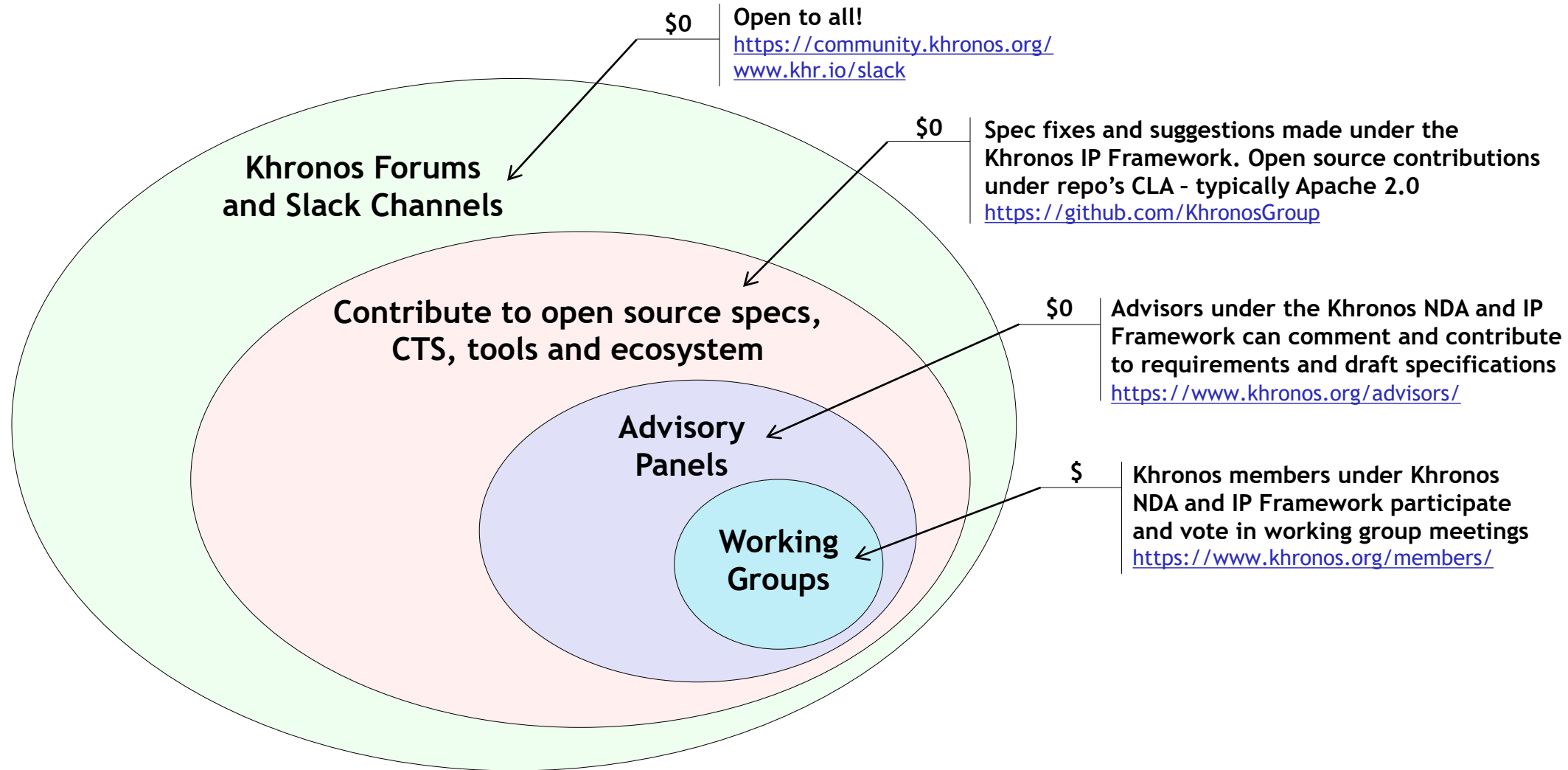**Silicon**

Open Source
Conformance Tests and
Adopters Programs

Non-profit organization -
Membership and Adopters
fees cover expenses

Invest where strong industry
momentum and relevance – let
Darwinism rule!

# Khronos Ecosystem Engagement

$0    **Open to all!**
      **https://community.khronos.org/**
      **www.khr.io/slack**

$0    **Spec fixes and suggestions made under the**
      **Khronos IP Framework. Open source contributions**
      **under repo's CLA – typically Apache 2.0**
      **https://github.com/KhronosGroup**

$0    **Advisors under the Khronos NDA and IP**
      **Framework can comment and contribute**
      **to requirements and draft specifications**
      **https://www.khronos.org/advisors/**

$     **Khronos members under Khronos**
      **NDA and IP Framework participate**
      **and vote in working group meetings**
      **https://www.khronos.org/members/**

**Khronos Forums
and Slack Channels**

**Contribute to open source specs,
CTS, tools and ecosystem**

**Advisory
Panels**

**Working
Groups**

# Benefits of Khronos membership

Gain early insights into industry trends and directions

Influence the design and direction of key open standards that will drive your business

Accelerate your time-to-market with early access to specification drafts

**Gather industry requirements for future open standards** → **Draft Specifications Confidential to Khronos members** → **Publicly Release Specifications and Conformance Tests**

Network with domain experts from diverse companies in your industry

State-of-the-art IP Framework protects your Intellectual Property

Enhance your company reputation as an industry leader through Khronos participation

# Thank You and Resources

- **Khronos is creating cutting-edge royalty-free open standards**
  - For 3D, compute, inferencing gaming

- **These slides and information on Khronos Standards**
  - www.khronos.org

- **Any company is welcome to join Khronos**
  - https://www.khronos.org/members/
  - We warmly welcome members from Australia and Asia

- **Dedicated developer resources**
  - Khronos Developer Forum: https://community.khronos.org/
  - Khronos Developer Slack Channel: www.khr.io/slack

- **We are happy to help answer any questions!**
  - Neil Trevett, Khronos President: ntrevett@nvidia.com, @neilt3d
  - Khronos Developer Relations, Kris Rose: kris@khronos.org, @kristoferrose

**WeChat: neiltrevett**