KHR SNOS®

Current and Planned Standards for Computer Vision and Machine Learning

010101010101010101010

810101010101010

Neil Trevett Khronos President NVIDIA VP Developer Ecosystems <u>ntrevett@nvidia.com</u> | <u>@neilt3d</u>

> embedded VISION



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

HUAWEI Imagination **EPIC** games (intel) QUALCOMM' SONY Veri Silicon AMDA Grm Google VALVE SAMSUNG De 3D Incorporated Contraction BARE BINOMIAL BURN (BOEING BRENWILL & BROADCOM CADENCE CALCT CAPASITY CODE CODE DE COD DisplayLink DREAMVIEW COMPENSION (A) ETTRE COMPENSION (A) STREET (pluto 🌾 R A Z E R 🐣 Red Hat RENESAS 서울대 학교 ⊕兆芯 🗿 shopify 🌣 SimplyAugmented Sing Socionext STREAM SynOPSYS* 🗱 University of 🗿 university of 🐺 University UX3) VALJ 🖬 💱 steen 🗇 where 🐨 swayfair 😢 XILINX. 💓 🔿 zSpace 📰

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, Vision Processing and Machine Learning

Khronos Active Initiatives

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







C LLADA...





Portable XR

Parallel Computation Vision, Inferencing, Machine Learning







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 Compute Acceleration Standards



This work is licensed under a Creative Commons Attribution 4.0 International License

SYCL Single Source C++ Parallel Programming



This work is licensed under a Creative Commons Attribution 4.0 International License

2

I

SYCL Implementations



ISO

This work is licensed under a Creative Commons Attribution 4.0 International License

S O Z

2

Т

 \mathbf{Y}



This work is licensed under a Creative Commons Attribution 4.0 International License

2

Т

 \mathbf{Y}

© The Khronos® Group Inc. 2019 - Page 7

Neural Network Exchange Formats

Before - Training and Inferencing Fragmentation

After - NN Training and Inferencing Interoperability



Two Neural Network Exchange Format Initiatives

NNEF	UNNX	ONNX and NNEF
Embedded Inferencing Import	Training Interchange	are Complementary ONNX moves quickly to track authoring framework updates NNEF provides a stable bridge from training into edge inferencing engines
Defined Specification	Open Source Project	
Multi-company Governance at Khronos	Initiated by Facebook & Microsoft	
Stability for hardware deployment	Software stack flexibility	

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



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



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

OpenVX Cross-Vendor Inferencing

OpenVX

S O Q Z

2

Т

 $\mathbf{\mathbf{\Sigma}}$

High-level graph-based abstraction for portable, efficient vision processing Graph can contain vision processing and NN nodes - enables global optimizations Optimized OpenVX drivers created and shipped by processor vendors Implementable on almost any hardware or processor with performance portability Run-time graph execution need very little host CPU interaction



Performance comparable to hand-optimized, non-portable code Real, complex applications on real, complex hardware Much lower development effort than hand-optimized code



This work is licensed under a Creative Commons Attribution 4.0 International License

© The Khronos[®] Group Inc. 2019 - Page 11

OpenVX Accelerated 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

S O Z V

H R

 $\mathbf{\Sigma}$

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



OpenVX/OpenCL Interop

OpenVX 1.3 Released October 2019



Open Source Conformance Test Suite https://github.com/KhronosGroup/OpenVX-cts/tree/openvx 1.3

ັທ

Z

Q

Т

 $\mathbf{\mathbf{Y}}$

Open Source OpenVX Tutorial and Code Samples

https://github.com/rgiduthuri/openvx_tutorial

Functionality Consolidation into Core 1.3

Neural Net Extension, NNEF Kernel Import, Safety Critical etc.

Deployment Flexibility through Feature Sets

Conformant Implementations ship one or more complete feature sets Enables market-focused Implementations - 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 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

OpenVX and OpenCV are Complementary

	OpenCV	SpenVX.
Implementation	Community driven open source library	Callable API implemented, optimized and shipped by hardware vendors
Scope	100s of imaging and vision functions Multiple camera APIs/interfaces	Tight focus on dozens of core hardware accelerated functions plus extensions and accelerated custom nodes. Uses external camera drivers
Conformance	Extensive OpenCV Test Suite but no formal Adopters program	Implementations must pass Khronos Conformance Test Suite to use trademark
IP Protection	None. Source code licensed under BSD. Some modules require royalties/licensing	Protected under Khronos IP Framework - Khronos members agree not to assert patents against API when used in Conformant implementations
Acceleration	OpenCV 3.0 Transparent API (or T-API) enables function offload to OpenCL devices	Implementation free to use any underlying API such as OpenCL. Can use OpenCL for Custom Nodes
Efficiency	OpenCV 4.0 G-API graph model for some filters, arithmetic/binary operations, and well-defined geometrical transformations	Graph-based execution of all Nodes. Optimizable computation and data transfer
Inferencing	Deep Neural Network module to construct networks from layers for forward pass computations only. Import from ONNX, TensorFlow, Torch, Caffe	Neural Network layers and operations represented directly in the OpenVX Graph. NNEF direct import, ONNX through NNEF convertor

This work is licensed under a Creative Commons Attribution 4.0 International License

HR

 $\mathbf{\mathbf{Y}}$

OpenVX Adoption for Inferencing



Source: Embedded Vision Alliance, Computer Vision Developer Survey, November 2019. © 2019 Embedded Vision Alliance



Rise in use of OpenVX for Inferencing

- Availability of NNEF Import and NN Extension
- Opportunity for silicon vendors to optimize inferencing solutions under a standard API
- Increasing deployment on embedded systems (in addition to HPC, Cloud, PCs, Mobile)

Trend will be accelerated by OpenVX 1.3's deployment flexibility

K H R O N O S

Primary Machine Learning Compilers



This work is licensed under a Creative Commons Attribution 4.0 International License

S O C Z

2

I



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

This work is licensed under a Creative Commons Attribution 4.0 International License

2

Т

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:

2

Т

 $\mathbf{\mathbf{\Sigma}}$

- 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 is Widely Deployed and Used



This work is licensed under a Creative Commons Attribution 4.0 International License

S O Z N

2

Т

 $\mathbf{\Sigma}$

OpenCL Evolution

OpenCL Extensions

Scratch-Pad Memory Management for DSPs Vulkan / OpenCL Interop Extended Subgroups etc.

Expanding Language Ecosystem

Tighter LLVM integration and cooperation Open source C++ for OpenCL Kernel Language SPIR-V 1.4 ingestion for compiler efficiency SPIR-V Extended debug info



OpenCL

May 2017 OpenCL 2.2 Improving Software Ecosystem Tool, libraries, ICD Loader Regular Maintenance Updates Spec clarifications, formatting and bug fixes <u>https://www.khronos.org/registry/OpenCL/</u>

Repeat Cycle for next Core Specification

Integration of Extensions plus New Core functionality



Target 2020

'OpenCL Next'

Focus for OpenCL Next is 'Deployment Flexibility'

- Flexible Profile so mobile and embedded vendors can ship customer-targeted functionality and be officially conformant
- Raise the bar for cross-vendor functionality on desktop and in cloud

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



This work is licensed under a Creative Commons Attribution 4.0 International License

2

Т

 $\mathbf{\Sigma}$

Need for New Camera Control API Standard?

- Khronos suspended work on OpenKCam standard several years ago
 - Mobile market went proprietary but embedded market has different needs
- OpenKCAM was aiming at advanced control of ISP and camera with cross-platform portability
 - Generate sophisticated image stream for advanced imaging & vision apps
 - Portable access to growing sensor diversity: e.g. depth sensors and sensor arrays
 - Cross sensor synch: e.g. synch of multiple camera and MEMS sensors
 - Advanced, high-frequency per-frame burst control of camera/sensor: e.g. ROI
 - Multiple input, output re-circulating streams with RAW, Bayer or YUV Processing



This work is licensed under a Creative Commons Attribution 4.0 International License

Т

Khronos Ecosystem Engagement



This work is licensed under a Creative Commons Attribution 4.0 International License

HR

 $\mathbf{\mathbf{\Sigma}}$

© The Khronos[®] Group Inc. 2019 - Page 23

Thank You and Resources

- Khronos is creating cutting-edge royalty-free open standards
 - For 3D, compute, vision, inferencing acceleration
- These slides and information on Khronos Standards
 - www.khronos.org

ຶ່

0° Z°

2

Т

- Any company is welcome to join Khronos
 - https://www.khronos.org/members/
- Dedicated developer resources
 - Khronos Developer Forum: <u>https://community.khronos.org/</u>
 - 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



