KHR SNOS®

APIs for Accelerating Vision and Inferencing An Overview of Options and Trade-offs

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



810101010101010

010101010101010**1010**



K H R S N O S S

Khronos is an open, non-profit, member-driven industry consortium developing royalty-free standards, and vibrant ecosystems, to harness the power of silicon acceleration for demanding graphics rendering and computationally intensive applications such as inferencing and vision processing

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

Some Khronos Standards Relevant to Embedded Vision and Inferencing



Embedded Vision and Inferencing Acceleration



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

 $\mathbf{\mathbf{Y}}$

Neural Network Training Acceleration



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

°S° O° Z

HR

 \mathbf{Y}

Embedded Vision and Inferencing Acceleration



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

 $\mathbf{\mathbf{Y}}$

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

Embedded Vision and Inferencing Acceleration



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

 $\mathbf{\mathbf{Y}}$

Primary Machine Learning Compilers



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

S O C Z

2

I

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



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



ົ

0°2 Z°

2

Т

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

Google's Multi-Level IR

Enables multiple domain-specific dialects within a common framework Experimenting with emitting to Vulkan/SPIR-V



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

°S° °° N°

2

Т

 \mathbf{Y}

© The Khronos[®] Group Inc. 2019 - Page 12

Embedded Vision and Inferencing Acceleration



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

 $\mathbf{\mathbf{Y}}$

PC/Mobile Platform Inferencing Stacks



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

HR

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

© The Khronos® Group Inc. 2019 - Page 14

Inferencing Engines



Platform Engines Android NNAPI

Microsoft WinML Apple CoreML



Desktop IHVs

AMD MIVisionX over MIOpen Intel OpenVINO over clDNN NVIDIA TensorRT over cuDNN



Acceleration APIs

Cross-platform Inferencing Engines





Both provide Inferencing AND Vision acceleration



HR

 \mathbf{Y}

Qualcomm Neural Processing SDK for Al



Mobile / Embedded

Arm Compute Library Huawei MACE Qualcomm Neural Processing SDK Synopsis MetaWare EV Dev Toolkit TI Deep Learning Library (TIDL) VeriSilicon Acuity



Almost all Embedded Inferencing Engines use OpenCL to access accelerator silicon

OpenVX and OpenCV are Complementary

	OpenCV	
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

OpenVX Cross-Vendor Inferencing

OpenVX

S O Q Z

2

Т

 $\mathbf{\Sigma}$

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



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

S O Q Z

H R

 $\mathbf{\Sigma}$

Kernel/Graph Import

- Provisional Extension
- Defines container for executable or IR code
- Enables arbitrary code to be inserted as a OpenVX Node in a graph



OpenVX/OpenCL Interop

OpenVX 1.3 Announced This Week!



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

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 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

Embedded Vision and Inferencing Acceleration



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

 $\mathbf{\mathbf{Y}}$

GPU and Heterogeneous Acceleration APIs



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

ັ້

O° Z°

2

Т

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

OpenCL is Widely Deployed and Used



Machine Learning Libraries

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

2

Т

 $\mathbf{\Sigma}$

SYCL Single Source C++ Parallel Programming

- SYCL 1.2.1 Adopters Program released in July 2018
 - https://www.khronos.org/news/press/khronos-releases-conformance-test-suite-for-sycl-1.2.1
- Multiple SYCL libraries for vision and inferencing
 - SYCL-BLAS, SYCL-DNN, SYCL-Eigen, SYCL Parallel STL



°S O N Z

2

I

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

E.g. complex ML frameworks can be

SYCL Implementations



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

Multiple Backend Support Coming

SYCL beginning to be supported on lowlevel APIs in addition to OpenCL e.g. Vulkan and CUDA <u>http://sycl.tech</u>

Intel's 'One API' Initiative uses SYCL

delivers-unified-programming-model-across-diverse-

architectures/#gs.bvdi6z

K H R S N S S S

Compiles C++-based SYCL

source files into code for both

CPU and a wide range of

compute accelerators

Con

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

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

© The Khronos[®] Group Inc. 2019 - Page 24

OpenCL Evolution



Focus for OpenCL Next is 'Deployment Flexibility'

Flexible Profile enables embedded vendors to ship targeted functionality for their customers and be officially conformant

Pervasive Vulkan

°S° O° Z°

2

Т

 \mathbf{Y}







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

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

Thank You and Resources

- Khronos Standards: OpenVX, NNEF, OpenCL, SPIR, SYCL, Vulkan and more...
 - Any company is welcome to join Khronos <u>www.khronos.org</u>
 - OpenVX 1.3: <u>www.khronos.org/openvx/</u>
 - SYCL: <u>http://sycl.tech</u>

ັ້

N

I

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

- Compilers for Machine Learning Conference Proceedings: <u>www.c4ml.org</u>
- MLIR: <u>www.blog.google/technology/ai/mlir-accelerating-ai-open-source-infrastructure/</u>
- Neil Trevett: ntrevett@nvidia.com | @neilt3d

K H R S N O S



Benefits of Khronos membership