KHR SNOS® GROUP

APIs for Accelerating Embedded Vision and Inferencing

Industry Overview of Options and Trade-offs Neil Trevett Khronos President NVIDIA VP Developer Ecosystems ntrevett@nvidia.com | @neilt3d



010101010101010101010

810101010101010

embeddedworld2020

Exhibition&Conference

K H R S N G R O U P

Connecting Software to Silicon

arm EPIC GAMES Google (intel) Imagination Over 150 members worldwide VALVE Veri Silicon SAMSUNG SONY Khronos is 20 years old this 💕 3D Incorporated 🚰 acer 🚺 🔕 Almotive 🦚 Collibada com Almalence 🖗 amazon.com 🔥 AGI 🔨 ARCANE Argonne 🛆 AUTODESK. AXISA BASE BINOMIAL AUTON (ABOEING BEENWILL BROADCOM CADENCE CAPEASITY CEVA © codeplay' Codeweavers C>O Collins Aerospace Ontinental's OCOOHOM COREAVE Coretronic GCRANK & CTRL-Lobs Traunhofer 🛞 التعميم المعالي المعا معالي معالي المعالي المعال معالي معالي المعالي MEDIATEK HICCOSOft Migenius Minute Antiputer Misuberatory Misuberatory model model model matrox NEC NIHON A R X E N T 🔛 🥺 The PEΔKHILLSGROUP ρluto 📥 Red Hat RENESAS 4書目 4 - 6 兆芯 G shopify O SimplyAugmented Strip Socionext STREAM SYNOPSYS TAKUMI 登 International Constants 🌵 Texas INSTRUMENTS 🛟 thinci 🎞 🖾 🕲 ThreeKit tobi TURBOS WID 🧐 unity 🛞 💵 ICT ultraleap 🕻 🛗 BRISTOL 🖞 LINVERSITY OF 🕅 University UX3) VALIA 🚺 Visteon' 🗇 MWARE' (VIS) 🕷 Wayfair' 🐔 XILINX. 💓 🔿 zSpace Institute

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

Open, non-profit, member-driven industry consortium creating advanced, royalty-free interoperability standards for 3D graphics, augmented and virtual reality, parallel programming, vision acceleration and machine learning

K H RON OS

Embedded Vision and Inferencing Acceleration



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

 $\mathbf{\mathbf{Y}}$

Khronos Active Initiatives

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



3D Assets Authoring and Delivery







Portable XR









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

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

© The Khronos[®] Group Inc. 2020 - Page 4

Khronos Compute Acceleration Standards



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 |
|-------------------------------------|-----------------------------------|
| 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 |

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



SYCL Single Source C++ Parallel Programming



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

2

I

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



S O Z

2

I

 \mathbf{Y}

OpenVX Cross-Vendor Inferencing

OpenVX

S O Z V

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. 2020 - Page 12

OpenVX 1.3 Released October 2019



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 Conformance Test Suite

https://github.com/KhronosGroup/OpenVX-cts/tree/openvx_1.3

OpenVX Open Source Implementation & Apps

Open Source OpenVX Tutorial and Code Samples

https://github.com/rgiduthuri/openvx_tutorial https://github.com/KhronosGroup/openvx-samples



Open Source OpenVX 1.3 for Raspberry Pi

Raspberry Pi 3 Model B with Raspbian OS Memory access optimization via tiling/chaining Highly optimized kernels on 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 | |
|----------------|---|--|
| Implementation | Community driven open source library | Callable API implemented, optimized and shipped by hardware vendors |
| Conformance | Extensive OpenCV Test Suite but no formal Adopters program | Implementations must pass Khronos Conformance Test Suite to use trademark |
| 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 |
| 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 |
| 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 |
| 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 |

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

°S° O° Z°

HR

 $\mathbf{\mathbf{Y}}$

Primary Machine Learning Compilers



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

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

© The Khronos[®] Group Inc. 2020 - Page 17

OpenCL - Low-level Parallel Programing

- Low-level programming of heterogeneous parallel compute resources
 - One code tree can be executed on CPUs, GPUs, DSPs, FPGAs, Tensor Processors ...
- 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:

S O Q Z

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



OpenCL Evolution

OpenCL Extensions

Asynchronous Copies for DSPs Extended Subgroups Extended Versioning

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

Integration of successful Extensions plus new Core functionality



Focus on 'Deployment Flexibility' to reach more processors and platforms

OpenCL

Target 2020 OpenCL Next

Extension Pipeline

Vulkan/OpenCL Interop Recordable Command buffers Unique Device Ids ML Primitives Device Topology Unified Shared Memory



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

Repeat The Cycle

K H R S N O S

Engaging with the Khronos Ecosystem



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

HR

 \mathbf{Y}



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

2

Т

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

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

 $\mathbf{\Sigma}$

Thank You!

- Khronos is creating cutting-edge royalty-free open standards
 - For 3D, compute, vision, inferencing acceleration
- Information on Khronos Standards: <u>www.khronos.org</u>
- Any company is welcome to join Khronos: https://www.khronos.org/members/
- Neil Trevett: <u>ntrevett@nvidia.com</u> | <u>@neilt3d</u>



Benefits of Khronos membership