

This webinar will begin shortly at
19:00 CET | 13:00 EST | 10:00 PDT



Accelerating Machine Learning with OpenCL

May 11, 2022

How to Participate

Speaker Questions

During the presentations, please submit speaker questions using the Zoom Q&A button (not the chat button). At the end of the talk, our moderator will put as many questions as possible to the speakers

General Questions and Comments

Please use the Zoom Chat feature for logistical questions or if you are having issues with Zoom

Recording

We are recording this webinar and will publicly post the link at the [Forum Home Page](#)

Survey

To help us design future Khronos ML Forum events, we appreciate you completing the short survey form that we will distribute after the session

Accelerating Machine Learning with OpenCL



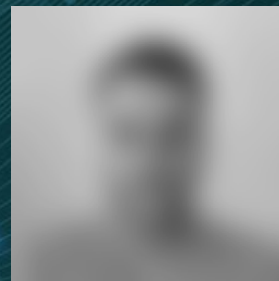
Khronos ML Forum
Neil Trevett, Khronos



A Case Study on OpenCL vs GPU Assembly
for Machine Learning Performance
Roy Oursler, Intel



Qualcomm Extensions for Advancing
Machine Learning Acceleration
Balaji Calidas, Qualcomm



Q&A Panel Moderator
Kevin Petit, Arm



Introduction to the Khronos Machine Learning Forum

Neil Trevett
Khronos President



Khronos Open Standards Mission



Open, royalty-free interoperability standards to harness the power of GPU, XR and multiprocessor hardware

3D graphics, augmented and virtual reality, parallel programming, inferencing and vision acceleration

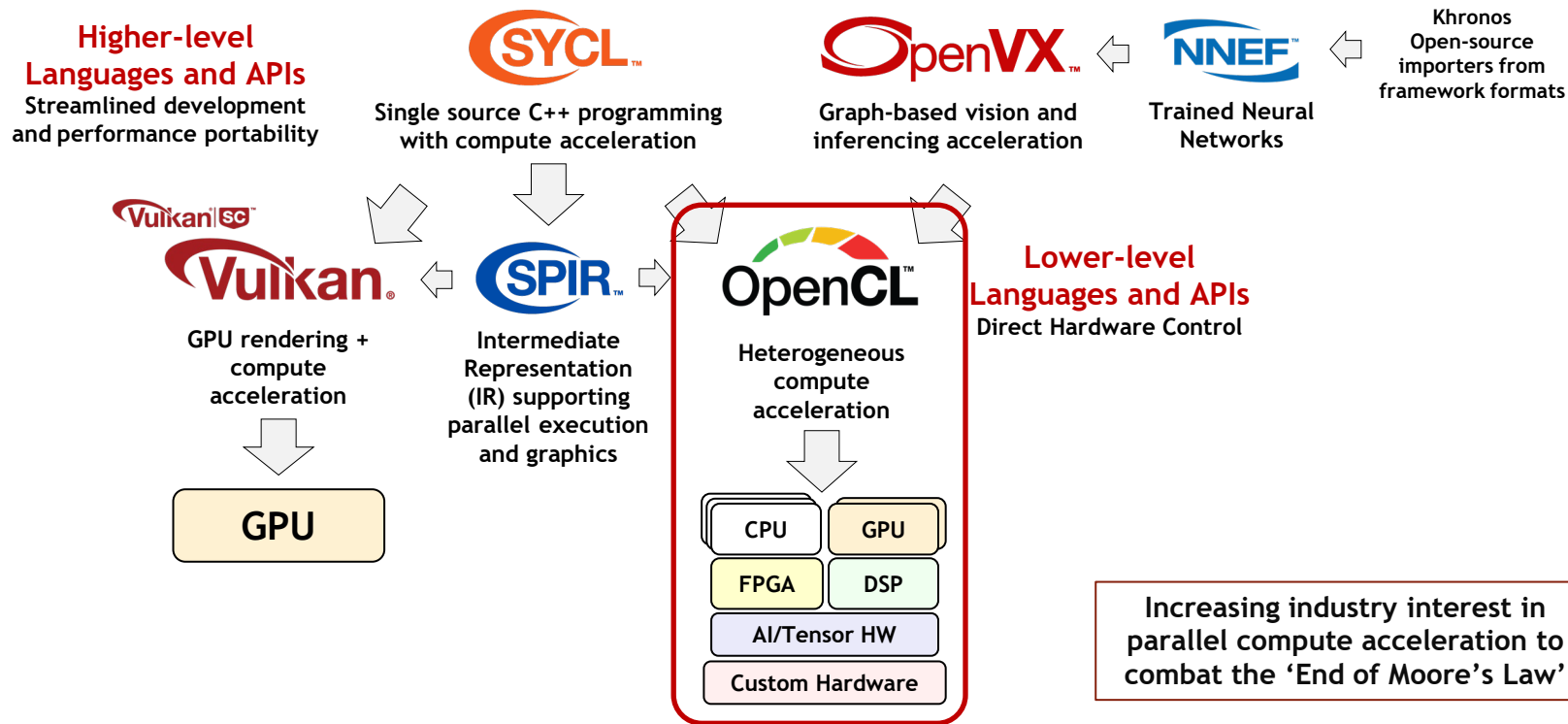
Non-profit, member-driven standards organization, open to any company

Proven multi-company governance and Intellectual Property Framework

Founded in 2000

~200 Members ~ 40% US, 30% Europe, 30% Asia

Khronos Compute Acceleration Standards



OpenCL and Machine Learning

Machine Learning
Compilers



Import Formats

Caffe, Keras,
MXNet, ONNX

TensorFlow Graph,
MXNet, PaddlePaddle,
Keras, ONNX

PyTorch, ONNX

TensorFlow Graph,
PyTorch, ONNX

Front-end / IR

NNVM / Relay IR

nGraph / Stripe IR

Glow Core / Glow IR

XLA HLO

Output

OpenCL, LLVM,
CUDA, Metal

OpenCL,
LLVM, CUDA

OpenCL
LLVM

LLVM, TPU IR, XLA IR
TensorFlow Lite / NNAPI
(inc. HW accel)



Common Steps

1.Import Trained
Network Description

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

3. Decompose to primitive
instructions and emit programs
for accelerated run-times

Machine Learning Compilers and Frameworks using OpenCL Acceleration

Inferencing Libraries and Frameworks

Alibaba MNN
Arm Compute Library
Baidu PaddlePaddle/Paddle-Lite
Berkeley Caffe
Intel c1DNN and OpenVINO

Google TensorFlow and NNAPI

SYCL-DNN
Synopsis MetaWare EV
Texas Instruments DL Library (TIDL)
VeriSilicon Acuity
Xiaomi Mace

Embedded NN Compilers

CEVA Deep Neural Network (CDNN)
Cadence Xtena
Neural Network Compiler (XNNC)



OpenCL Extension Pipeline

A significant percentage of OpenCL extensions are used for Machine Learning acceleration

New OpenCL KHR extensions shipped since IWOCL 2021

Subgroup rotate extension for efficient data exchange among work-items

Workgroup Uniform Arithmetic for new work-group scan and reduction operators

Command Buffers Record and Replay (Provisional)

Asynchronous DMA

Expect Assume Hints

Integer Dot Product

External memory objects and semaphores (Provisional)

EXT and Vendor Extension Pipeline

Generalized Image from buffer (EXT)

Unified Shared Memory (EXT)

Floating Point Atomics (EXT)

Cooperative Matrices (EXT)

Machine Learning Operations (Qualcomm)

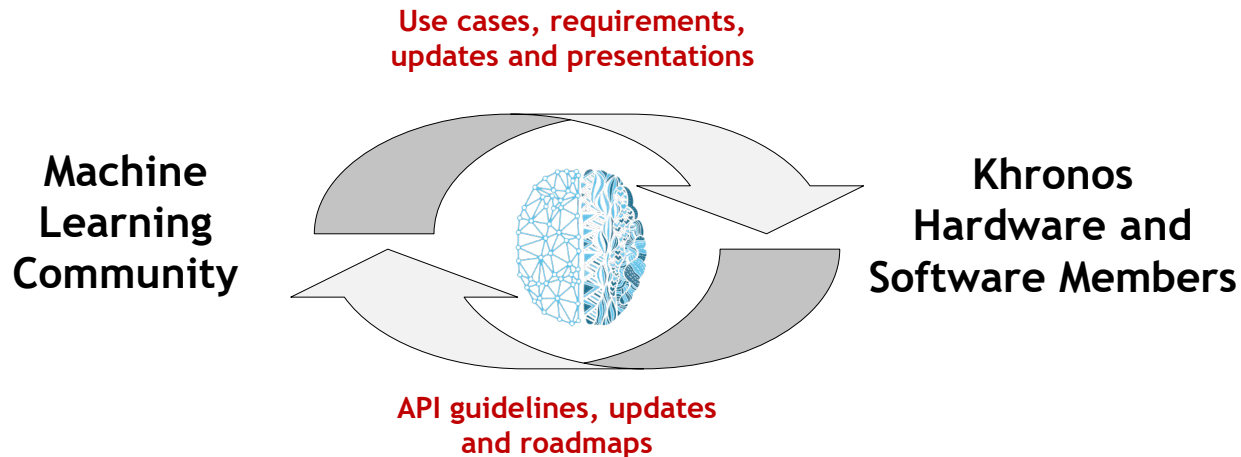
....

Khronos Machine Learning Forum

Productive ongoing communication and cooperation on ML Acceleration
... between Machine Learning *hardware* and *software* communities

Forum is free to join, no NDA or IP commitments

Dedicated meetings, email and slack channels for group communication



Machine Learning Forum Meeting Series



Forum Member Meetings will start in July 2022
Input and requests for specific topics welcome!

All the information you need to join!
<https://www.khronos.org/machine-learning>



Accelerating Machine Learning with OpenCL



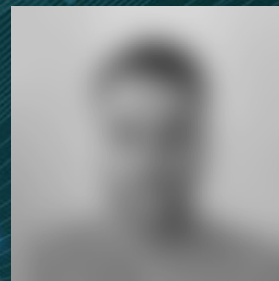
Khronos ML Forum
Neil Trevett, Khronos



A Case Study on OpenCL vs GPU Assembly
for Machine Learning Performance
Roy Oursler, Intel



Qualcomm Extensions for Advancing
Machine Learning Acceleration
Balaji Calidas, Qualcomm



Q&A Panel Moderator
Kevin Petit, Arm

Accelerating Machine Learning with OpenCL

11 May 2022

Ask the Experts

Use the Zoom Q&A feature to ask your questions



IWOCL & SYCLcon 2022

Accelerating Machine Learning with OpenCL

11 May 2022

Thank you!

<https://www.khronos.org/events/accelerating-machine-learning-with-opencl>