

04/04/2019

mcharleb@qti.qualcomm.com

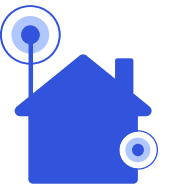
Bangkok, Thailand

Qualcomm

Inferencing at the Edge and Fragmentation Challenges

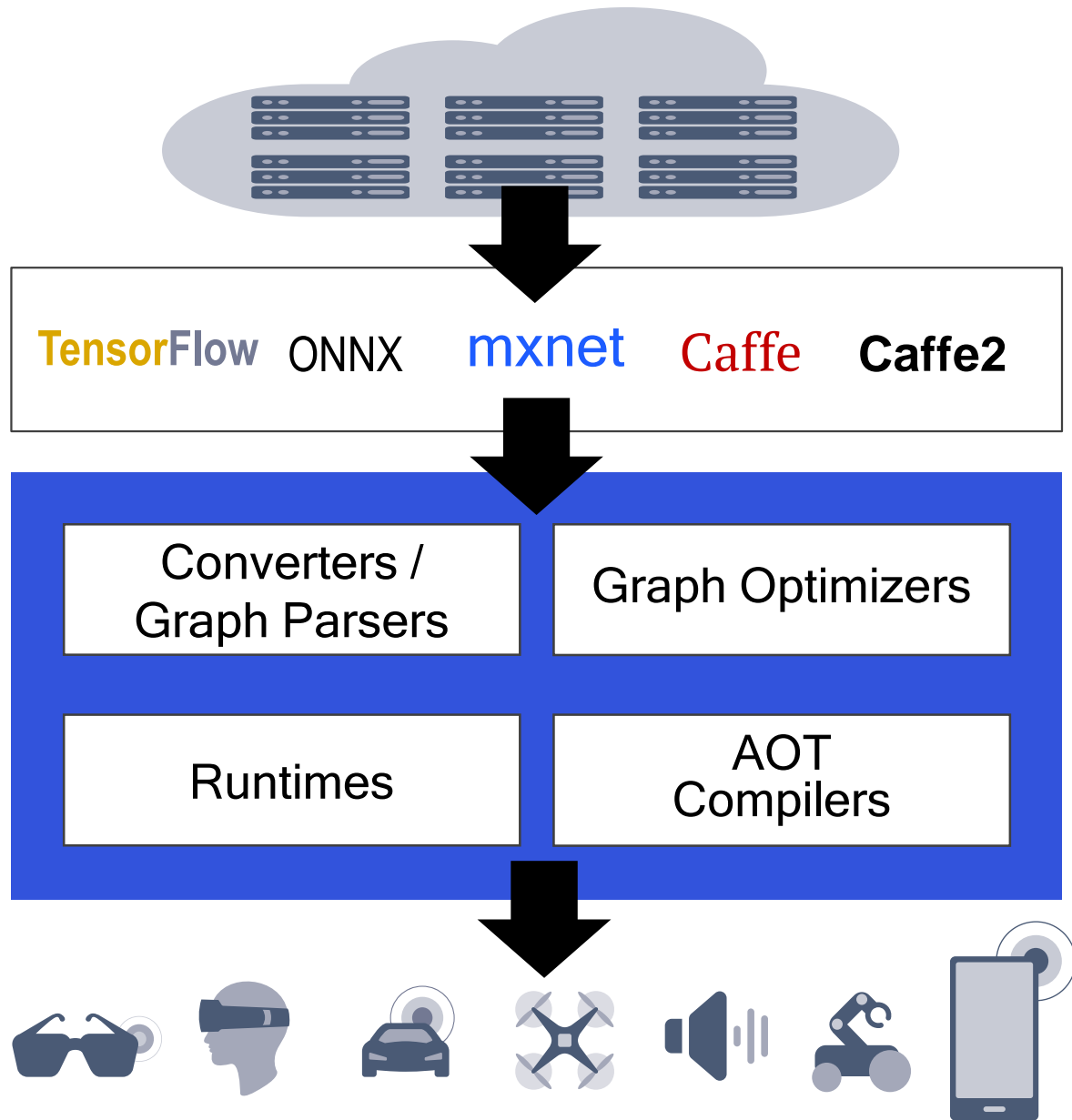
Mark Charlebois

Director Engineering
Qualcomm Technologies, Inc.



Agenda

- Overview of Inferencing at the Edge
 - Deployment Scenarios
 - Application Development Options
- Runtimes
- Graph Compilers
 - Four Levels of IR
- Leveraging Runtimes in Compiler Frameworks
- Conclusions

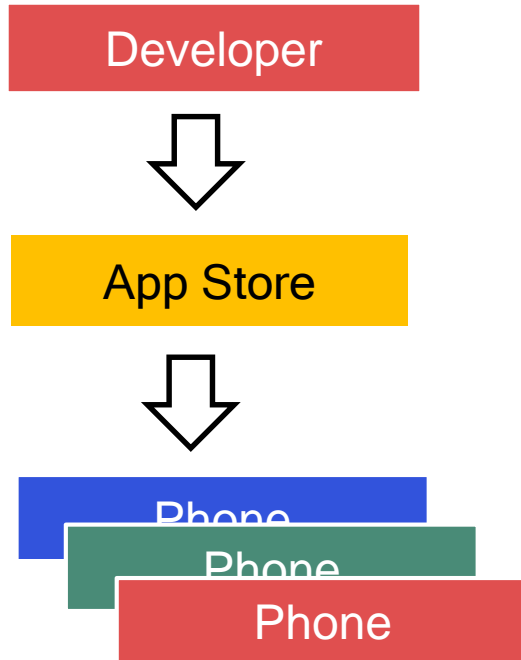


Inferencing at the Edge

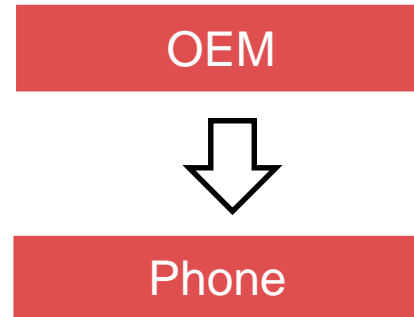
- Models developed on servers, deployed to edge devices
- Devices can range from MCUs to edge servers
- Simple networks can run on MCU
- Huge diversity of HW at the edge for inferencing acceleration
- Growing number of runtimes and graph compilers

Deployment Scenarios

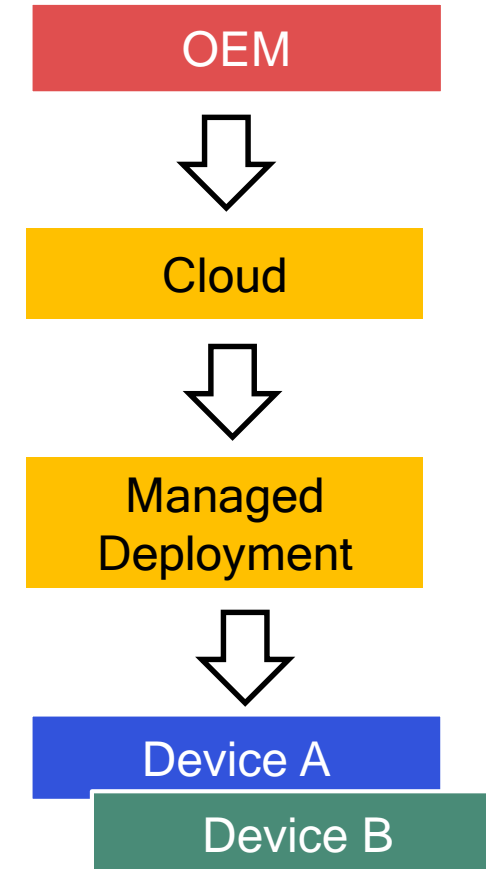
Known Models,
Unknown Devices



Known Models,
Known Device

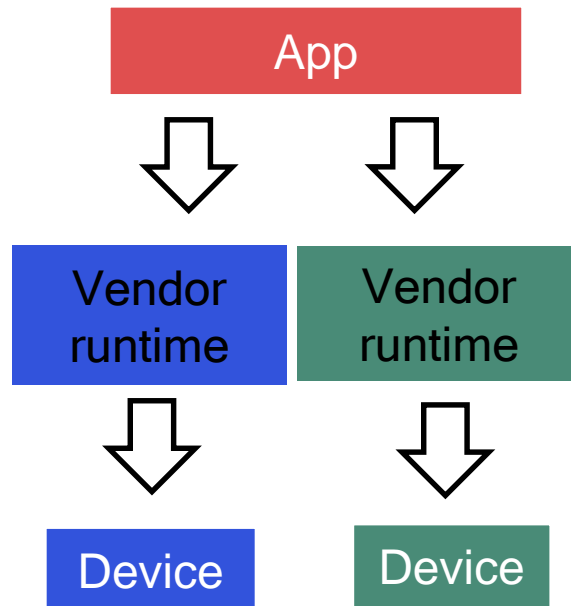


Known Models,
Known Devices

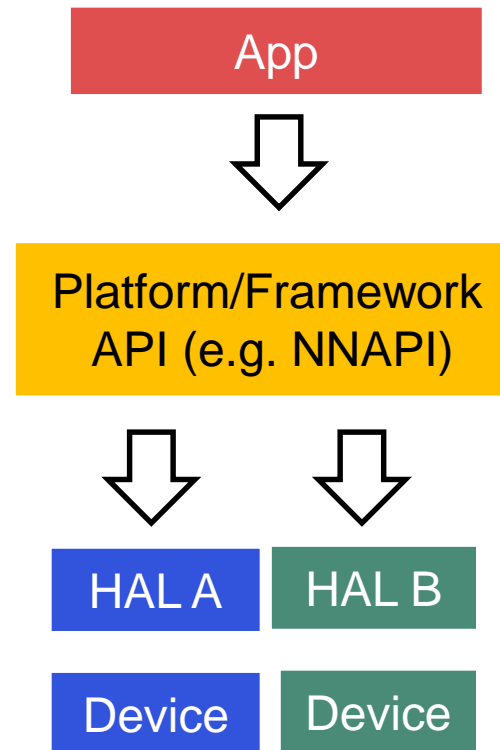


Application Development Options

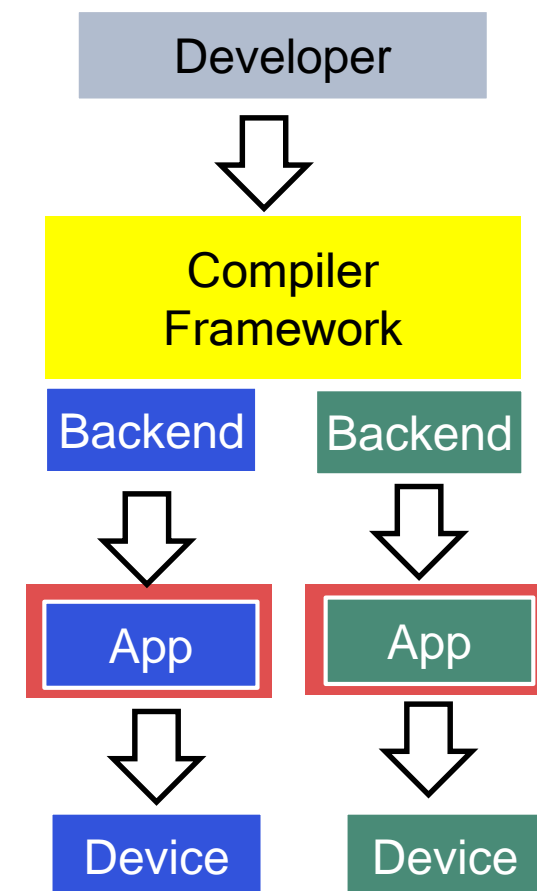
App Using Vendor SDKs



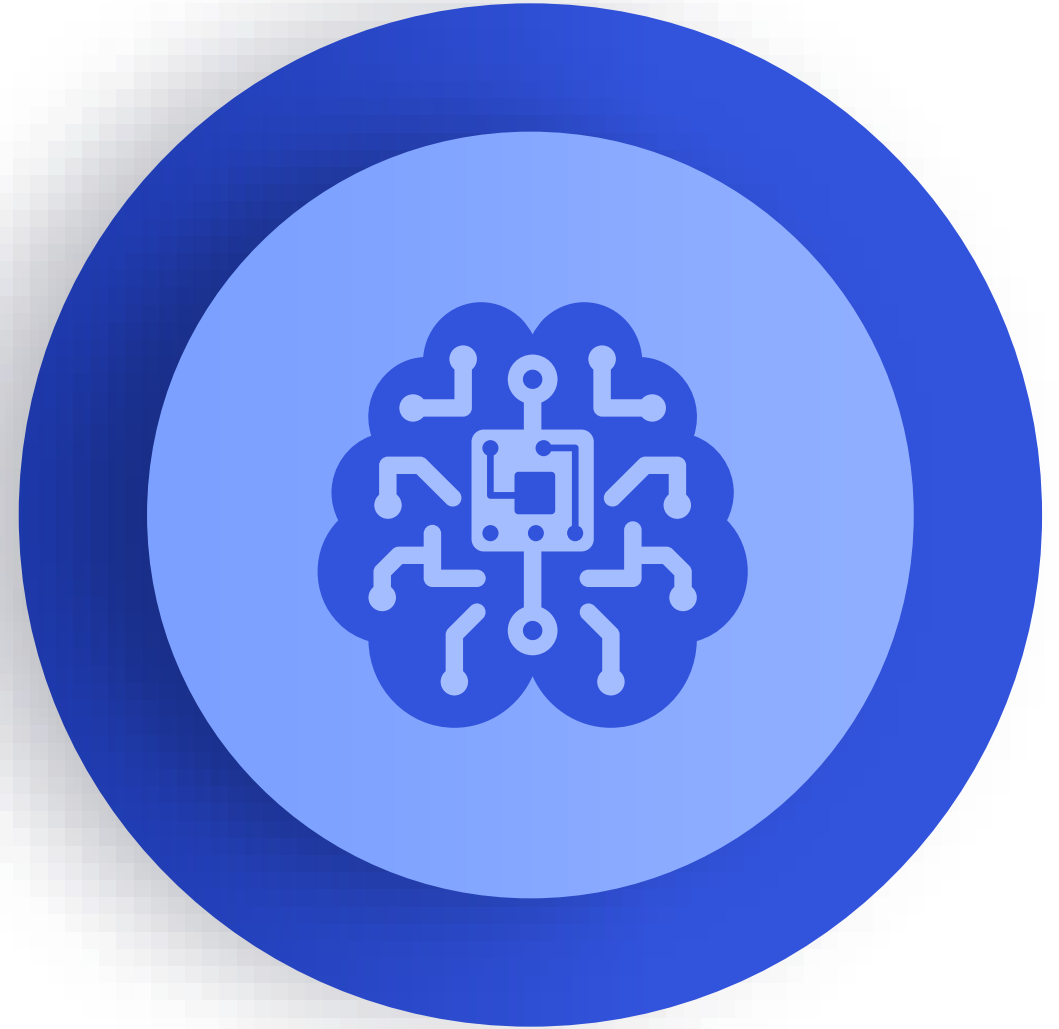
App using OS / Framework API



App Using Compiled Model



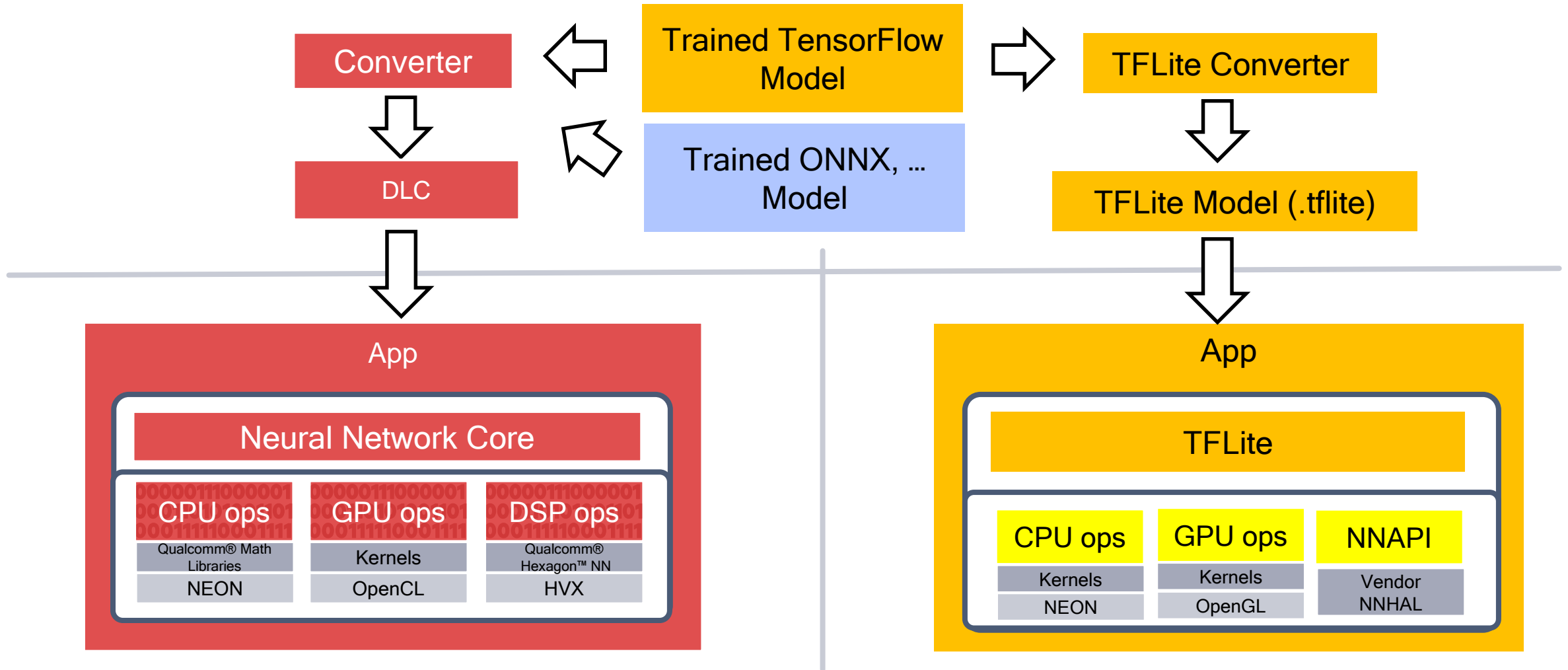
Runtimes



App Examples

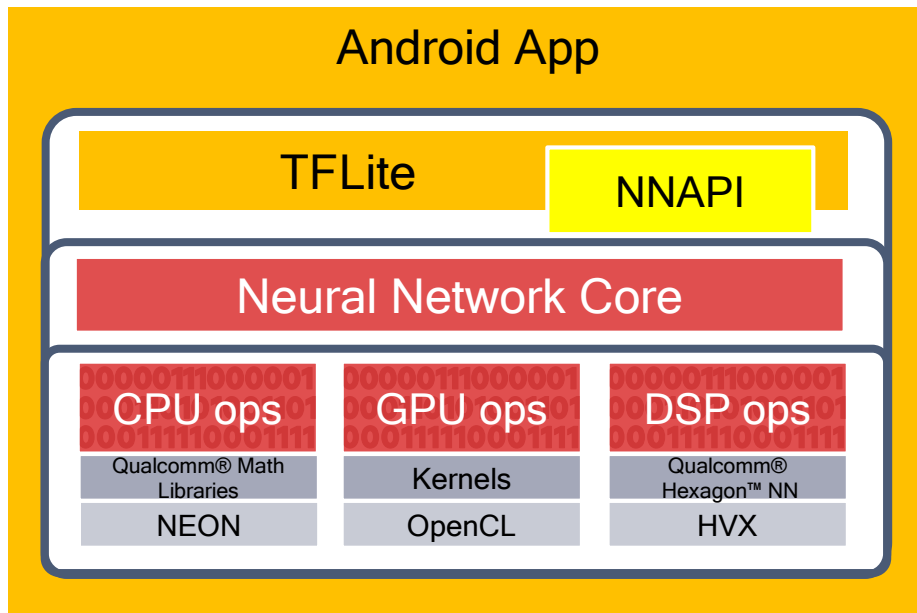
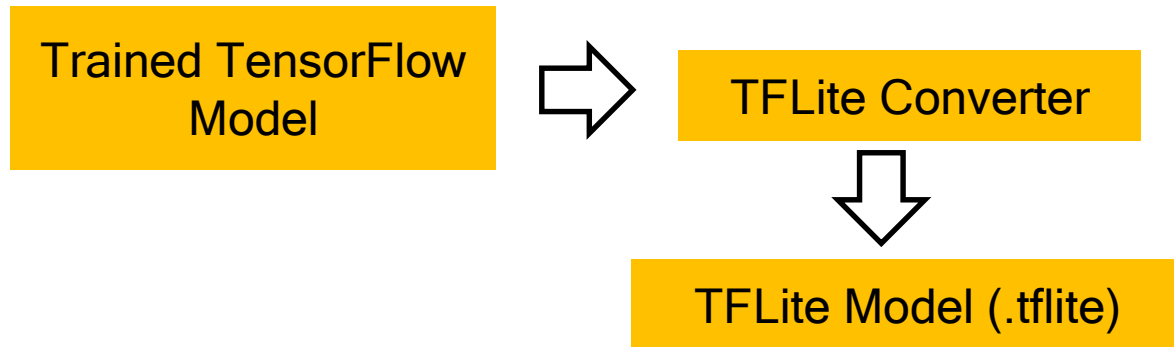
- Qualcomm® Neural Processing SDK

- TensorFlow Lite



Android App Example (NNAPI)

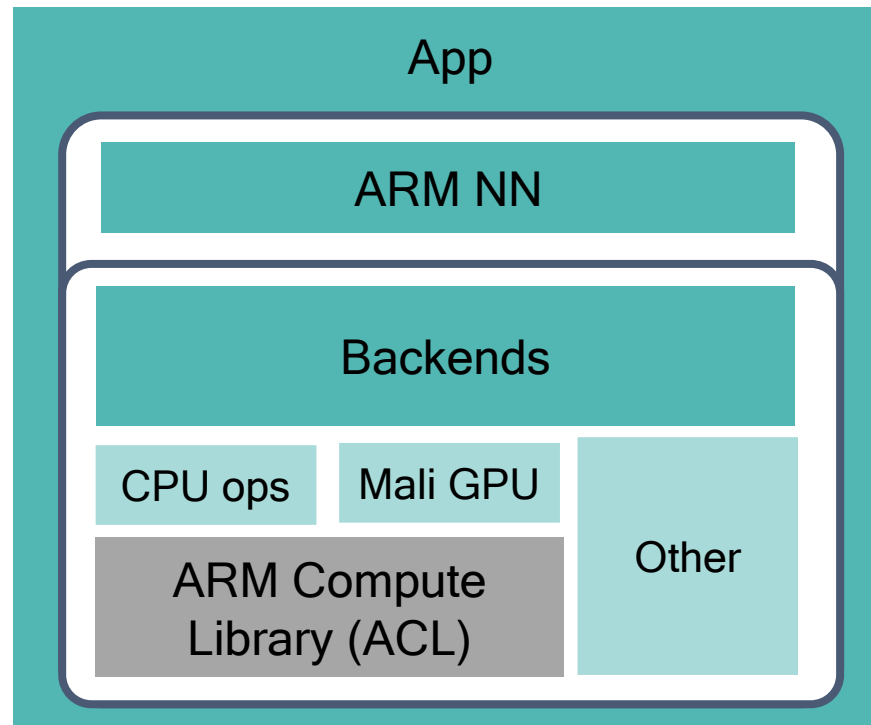
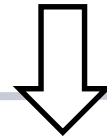
- TFLite App running on Qualcomm® device using NNAPI



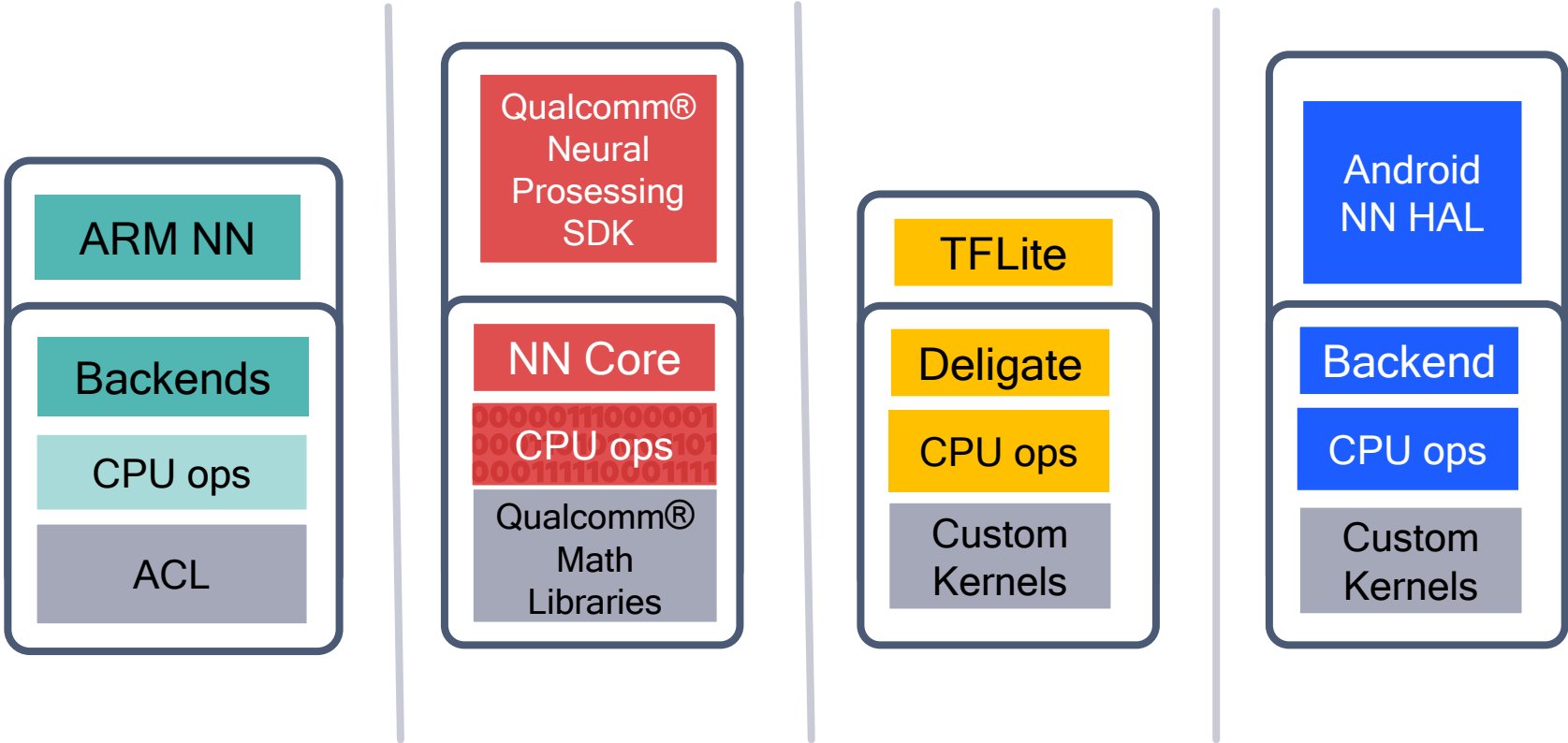
ARM NN Runtime Example

Trained TensorFlow Model

Trained ONNX, ... Model



CPU Runtime Fragmentation



Graph Compilers



Graph Compilers

ONNX Runtime

TVM

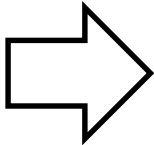
nGraph

ONNC

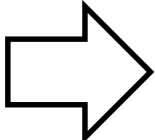
XLA

GLOW

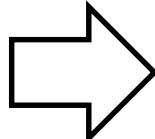
e.g. TensorFlow



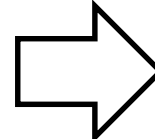
Graph Compiler Framework



CPU Binary

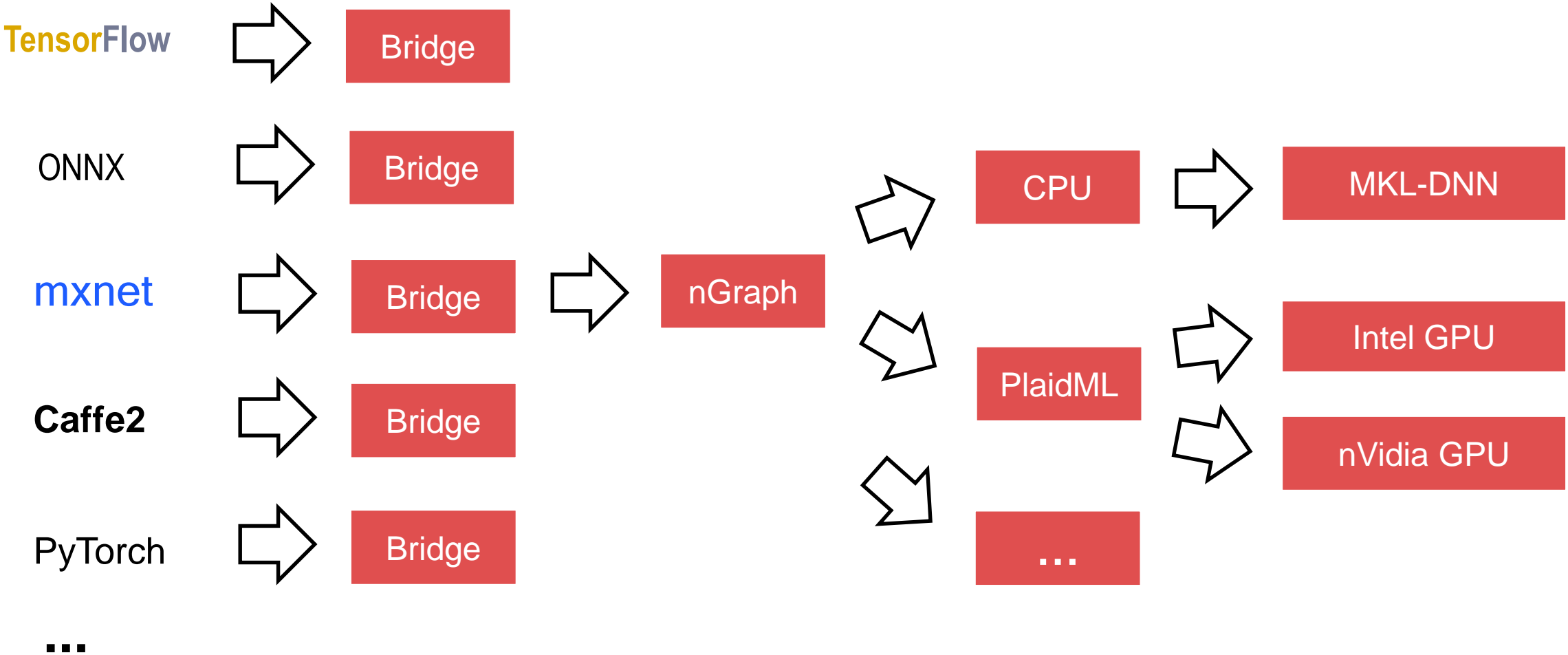


GPU Binary

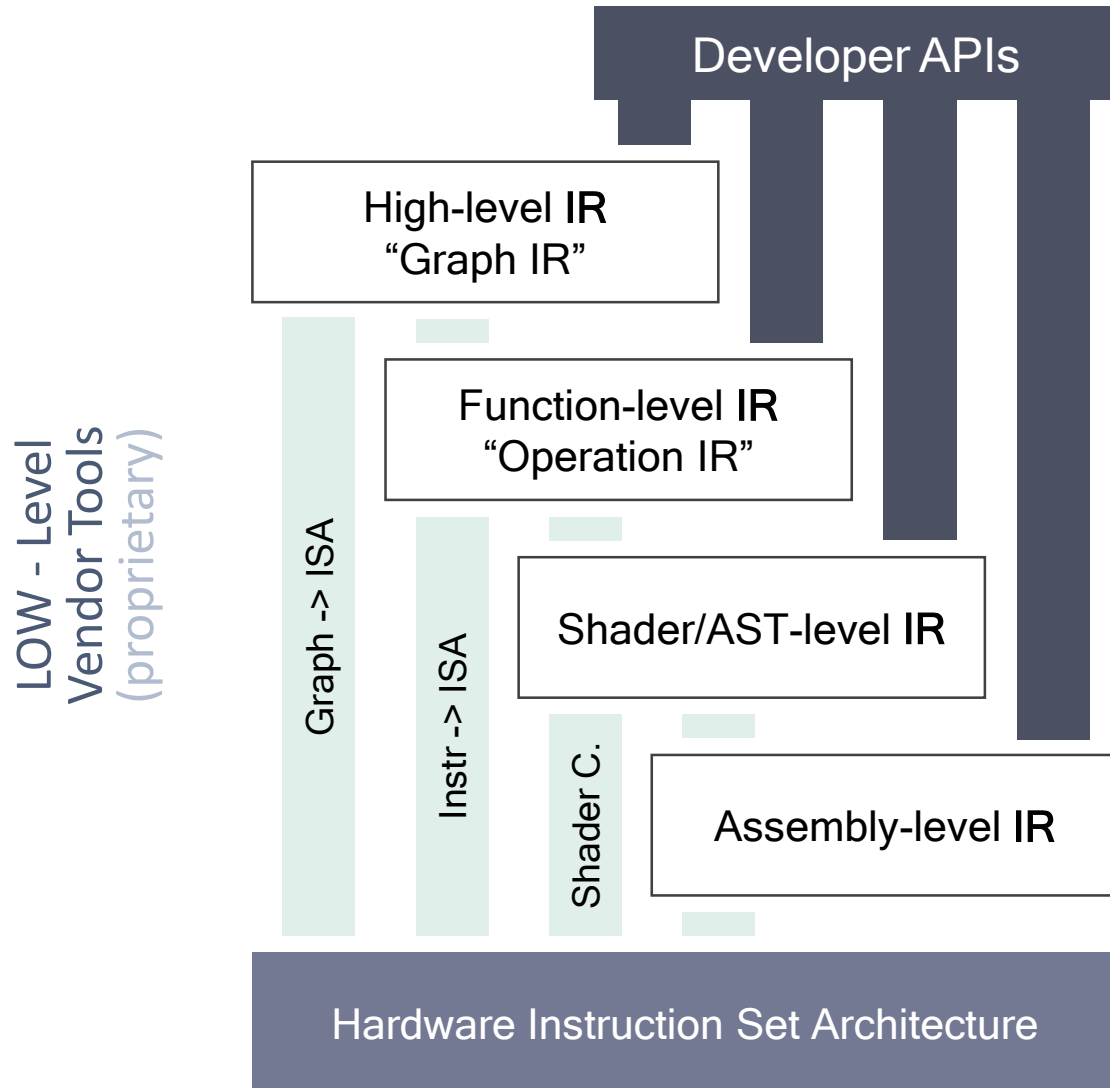


DSP Binary

Framework Layer Cake - Example nGraph



Four Levels of IR



- Parse Graph, first pass optimization and partitioning
- Graph segment optimization
- Memory reuse (non-backend specific)
- Layout transformation

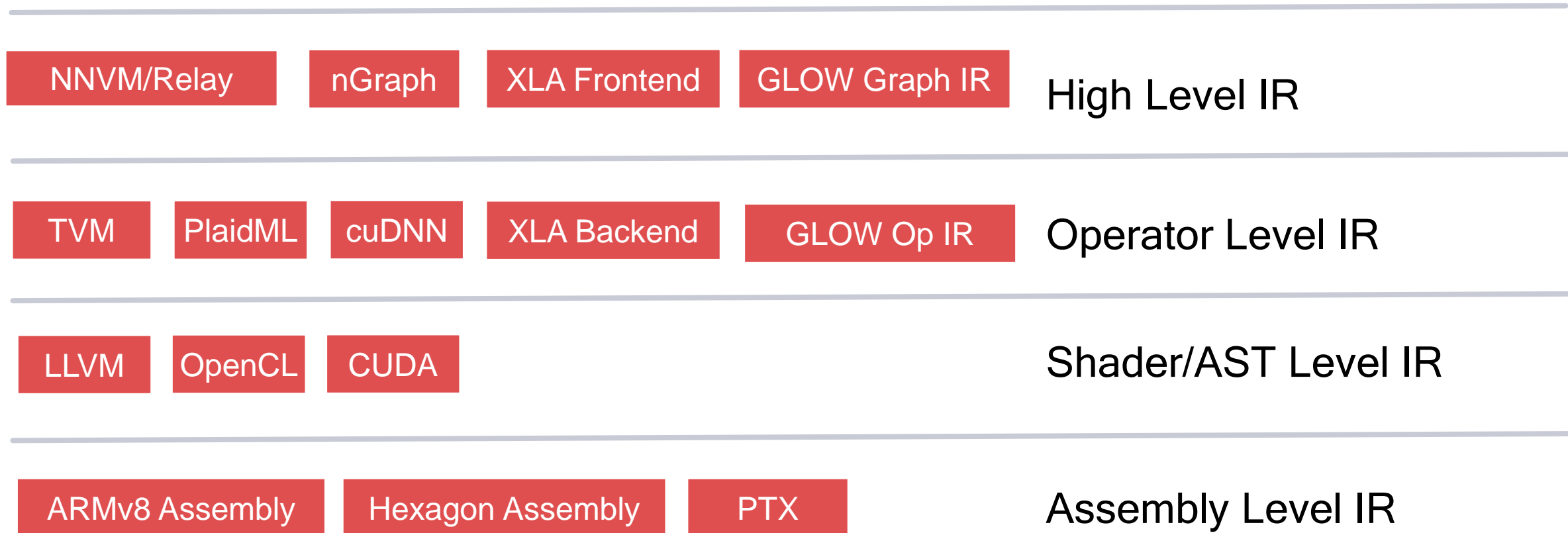
- HW specific subgraph optimizations
- HW optimizations:
 - Explicit memory layout, parallelization pattern

- Generate shaders, DSL code, bitcode
- Integrate hand optimized kernels

- Generate assembly and schedule for target ISA

Projects Mapped to Four Levels of IR

TensorFlow ONNX mxnet Caffe2 PyTorch ...



Leveraging Runtimes in Compiler Frameworks

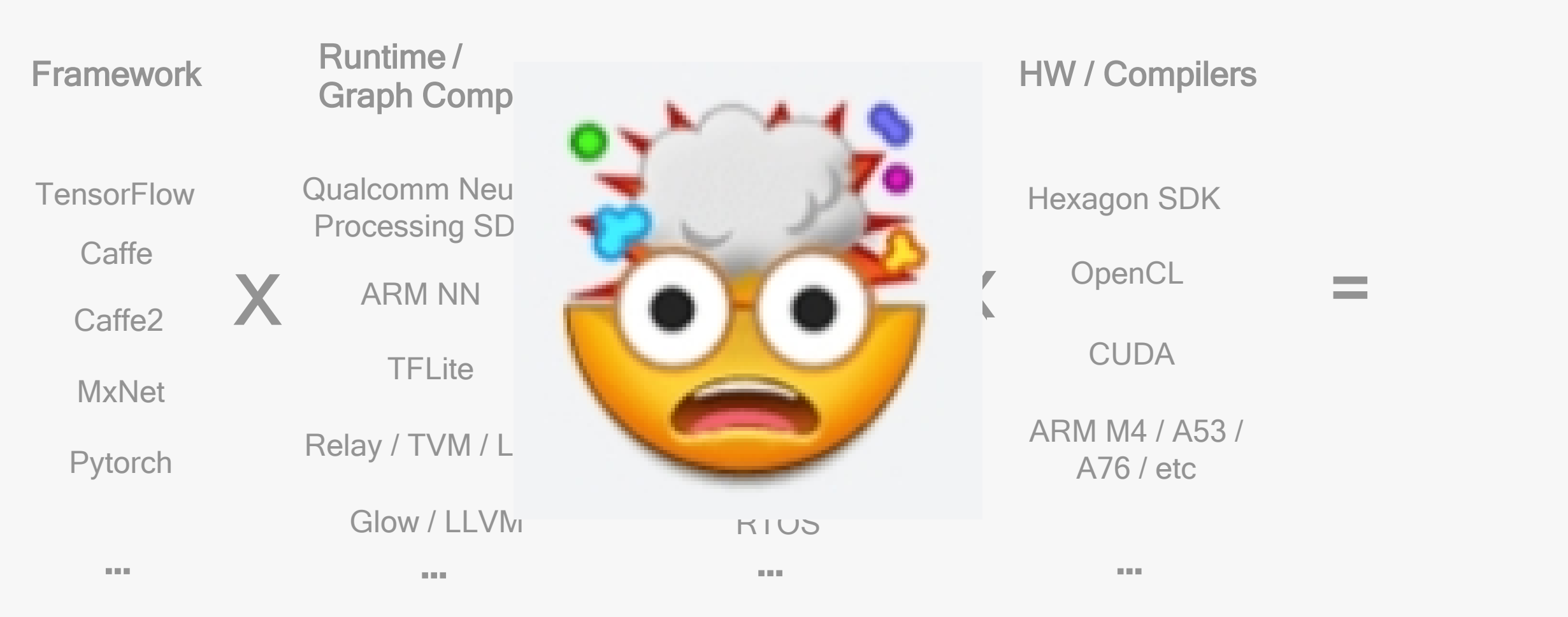


Frameworks, Platforms, HW, All Growing

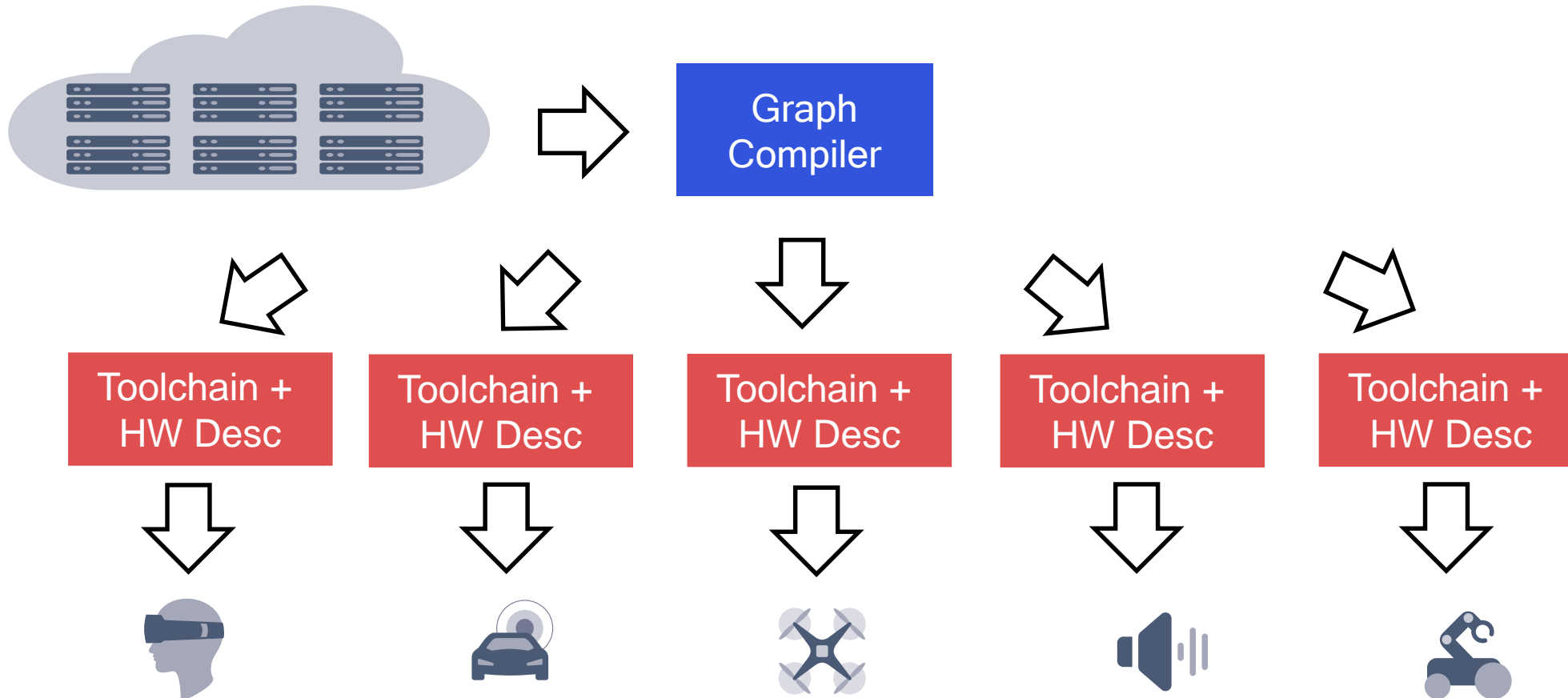
Framework	Runtime / Graph Compiler	OS	HW / Compilers
TensorFlow	Qualcomm Neural Processing SDK	Open Embedded	Hexagon SDK
Caffe	X	ARM NN	X
Caffe2		Android	
MxNet	TFLite	Ubuntu / Fedora / etc	OpenCL
Pytorch	Relay / TVM / LLVM	Windows	CUDA
	Glow / LLVM	RTOS	ARM M4 / A53 / A76 / etc
...

=

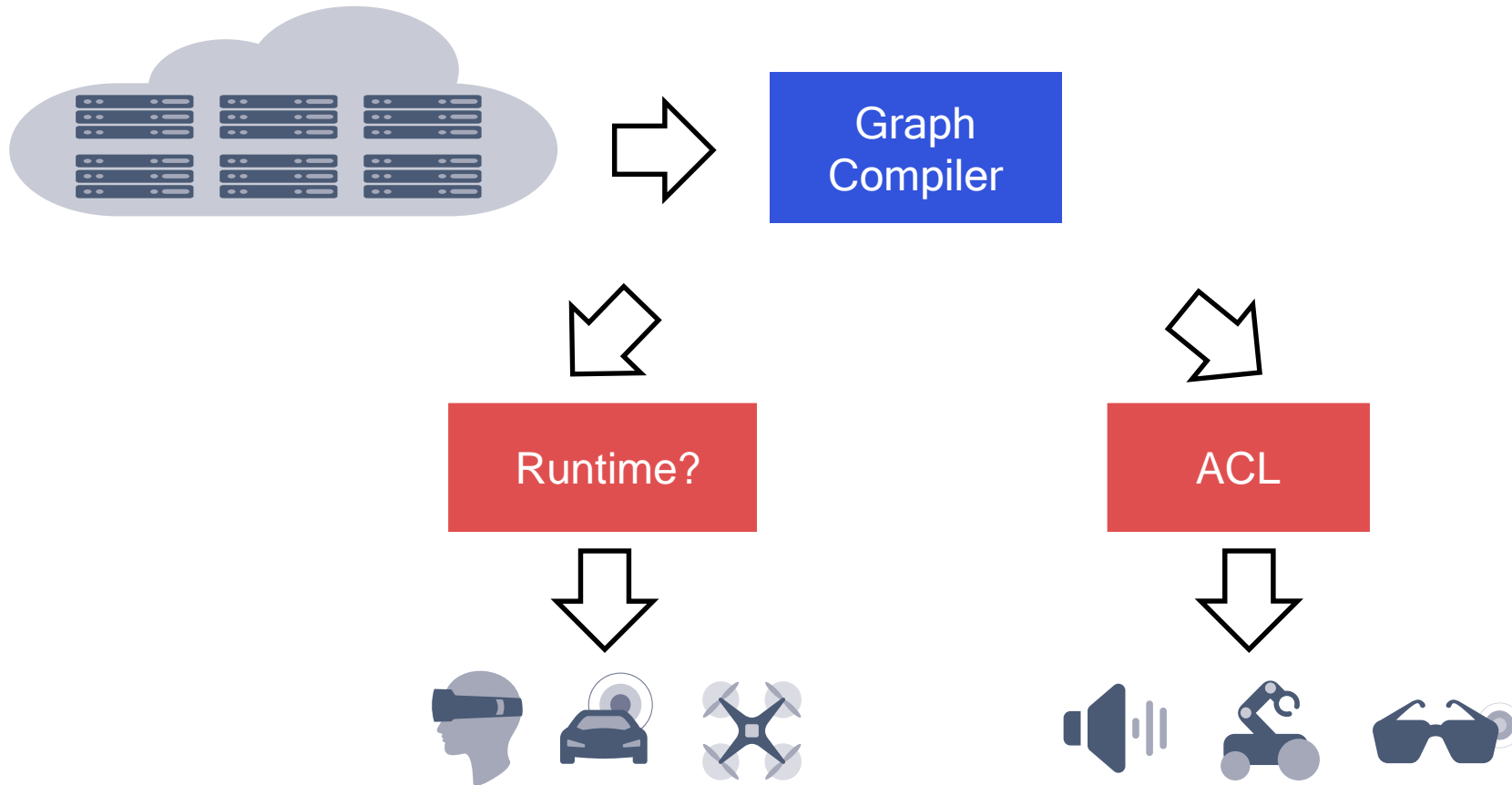
Frameworks, Platforms, HW, All Growing



HW Diversity Is a Scalability Issue



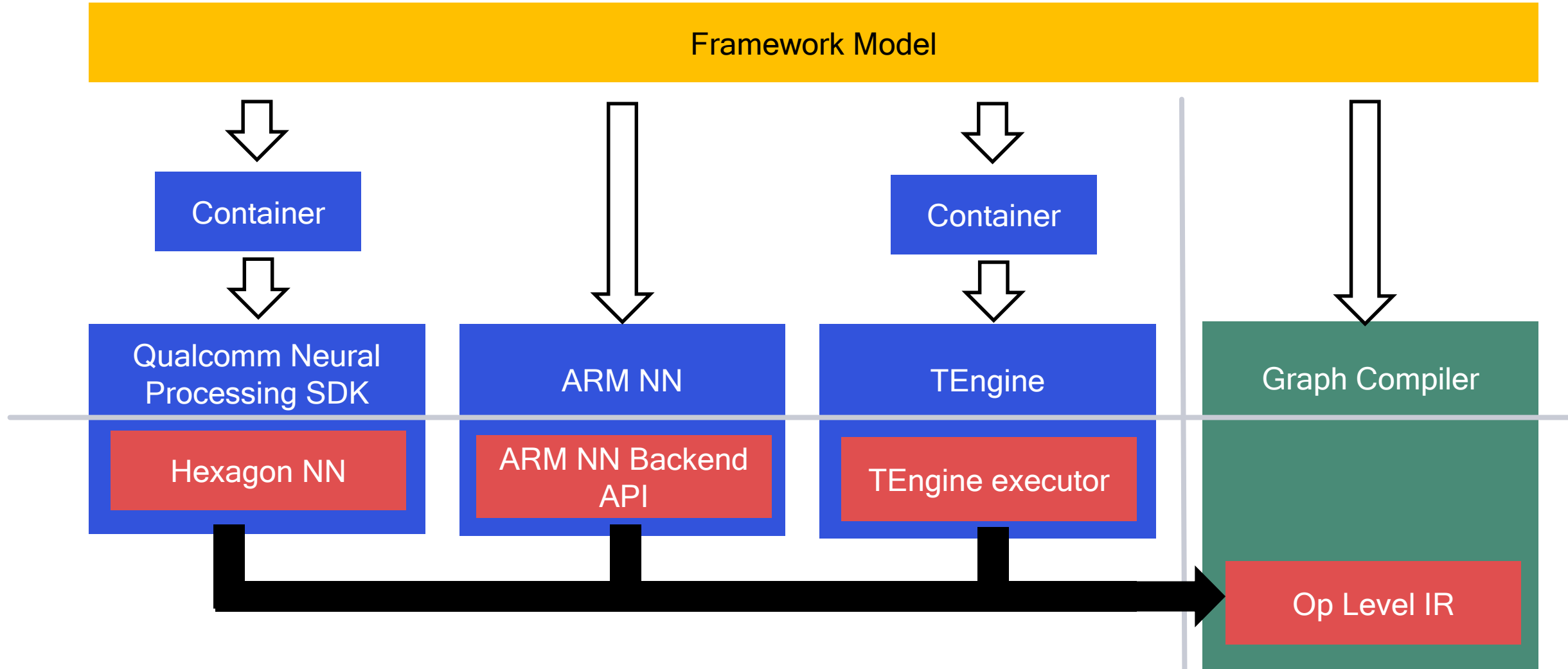
Tradeoff Some Performance for Scale



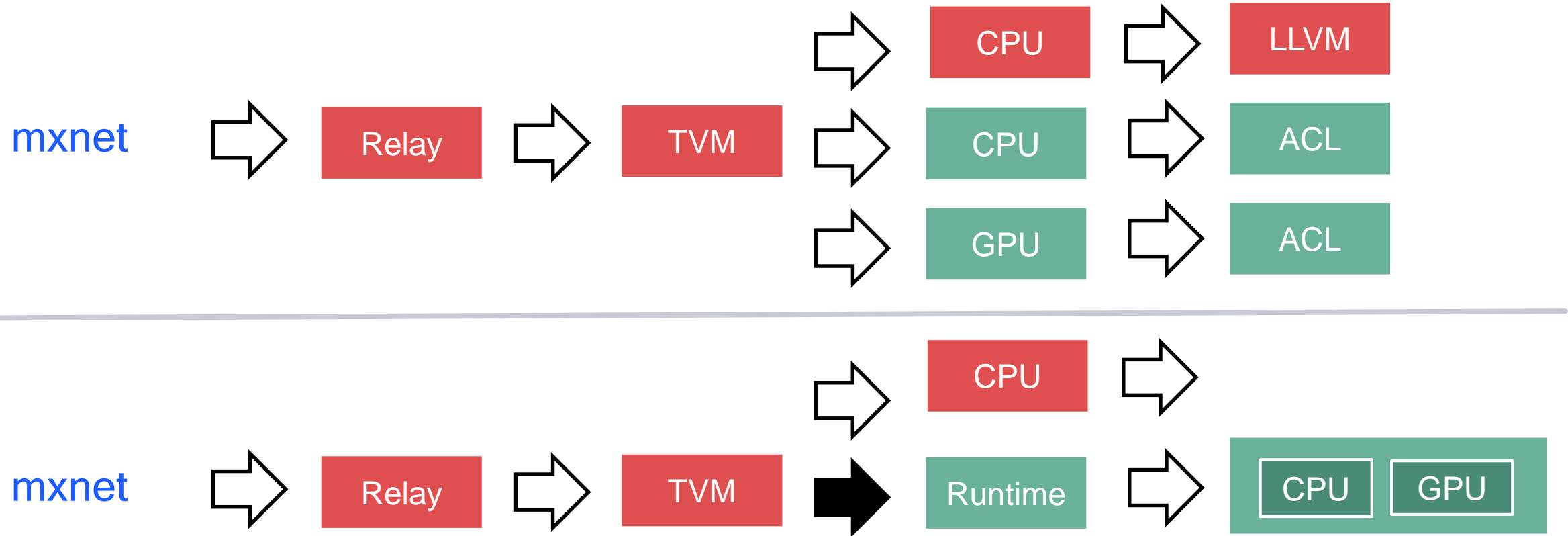
Conclusions



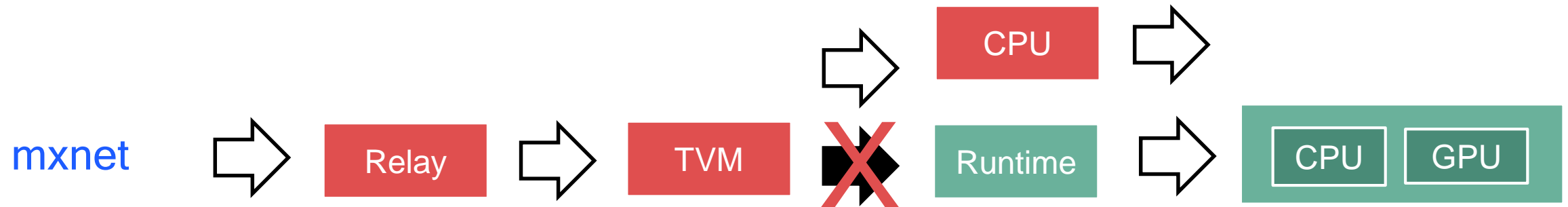
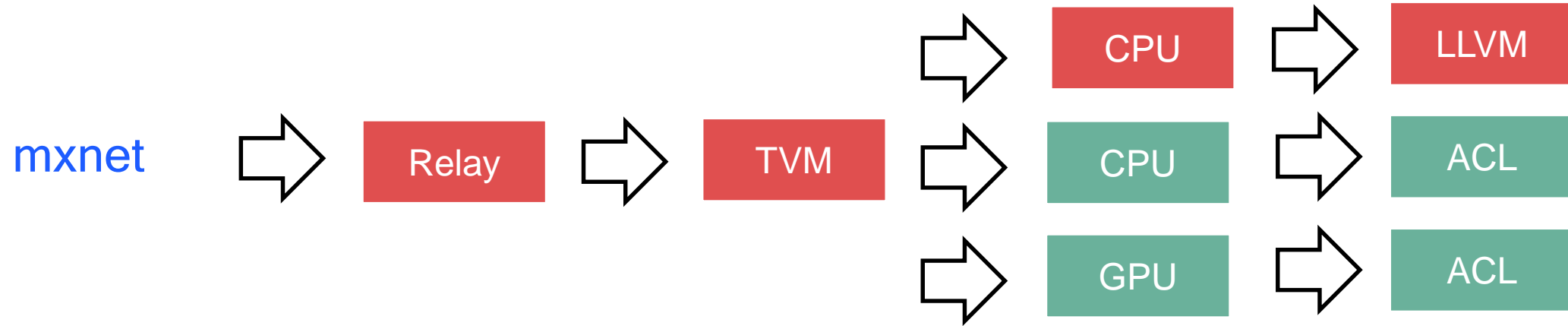
Runtime Components Map to Op Level IR



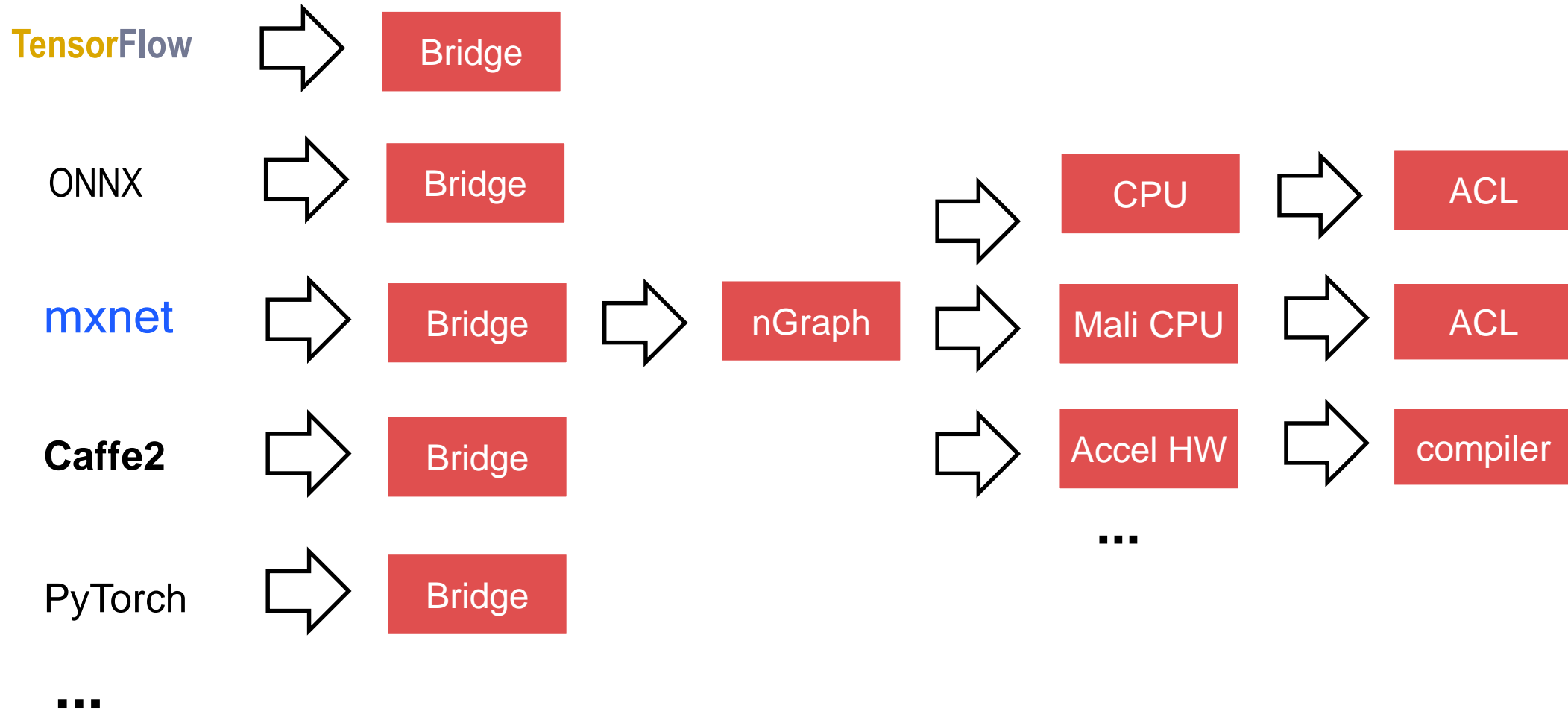
Using Full Runtime as a HW Backend



Using Full Runtime as a HW Backend



ARM Inferencing Edge Server



Addressing Graph Compiler IR Fragmentation




- Each compiler framework makes different tradeoffs, but eventually may be able to share components: e.g. code generation for OpenCL, LLVM, etc
- Too many frameworks, compilers and formats to address them all.
- Too much disparate HW for all frameworks to support.
- TensorFlow/TFLite and ONNX formats can provide the most scale for edge device inferencing runtimes.
- If Compiler Frameworks supported a common runtime backend API (like ARM NN Backend API) to bind to operator IR would enable graph compilers to support more edge devices with optimized backends, and would provide a common API for new backends (e.g. Hexagon NN) vs individual ports to each project.

Addressing CPU Runtime Fragmentation

- Make ACL the “best of breed” CPU Runtime
- Consolidate TFLite CPU runtime, Android NN CPU runtime, and ARM NN CPU backend
- Paves the way for others to follow: TEngine, MACE, Qualcomm® Neural Processing SDK, ...



Thank you!

Follow us on:   

For more information, visit us at:

www.qualcomm.com & www.qualcomm.com/blog

Nothing in these materials is an offer to sell any of the components or devices referenced herein.

©2019 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

Qualcomm is a trademark of Qualcomm Incorporated, registered in the United States and other countries. Qualcomm Hexagon is a product of Qualcomm Technologies, Inc. and/or its subsidiaries. Other products and brand names may be trademarks or registered trademarks of their respective owners.

References in this presentation to “Qualcomm” may mean Qualcomm Incorporated, Qualcomm Technologies, Inc., and/or other subsidiaries or business units within the Qualcomm corporate structure, as applicable. Qualcomm Incorporated includes Qualcomm’s licensing business, QTL, and the vast majority of its patent portfolio. Qualcomm Technologies, Inc., a wholly-owned subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of Qualcomm’s engineering, research and development functions, and substantially all of its product and services businesses, including its semiconductor business, QCT.