# embedded VISION SUMMIT 2018

# Neural Network Compiler:
# Enabling Rapid Deployment of DNNs on Low-Cost, Low-Power Processors

**cadence**®

Megha Daga

05/23/2018

- Industry Trend
- Software Challenges
- Software Solution:
    - Static Tools
    - Dynamic Tools
- Tensilica® Xtensa® Neural Network Compiler (XNNC): Static Tool Support
- End-to-End deployment example
- Android Neural Network API (ANN): Dynamic Tool Support
- Conclusion

*"Alexa, Add a 2pm meeting to my calendar"*

# Where is Processing Happening?

- AI in cloud
  - Adds latency
  - Undependable connectivity
  - Breaches privacy
  - High cost
- Trend: AI On-Device
  - Face detection
  - Seeing AI
  - Speech to Text

cādence®

- Requirements
  - Low Power
  - Low Area
  - Programmable
  - High Performance
  - Efficient Memory Management
  - Low Latency
  - Scalable

**AI Specialized Embedded Device**
**Example:**
*Tensilica® Vision P6,*
*Tensilica® Vision Q6,*
*Tensilica® Vision C5*

# Embedded Platform: NN Implementation Challenges

## Specialized Engines

- Optimum performance needs utilization of intrinsics

- Optimum performance needs most appropriate selection of ISA

## Market Needs

- Accelerate algorithm implementation time to keep up with the speed of product announcements

- Need for dynamic platform to adapt to rapidly changing neural network algorithms

## Optimization

- Need for optimization due to huge compute and bandwidth needs from neural networks

- Need for optimum memory management for bandwidth and power reduction
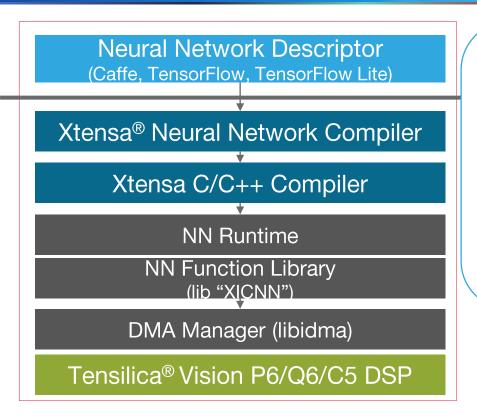
cādence®

## Offline

- Varied framework support
- Pre-defined NN models
- Need most optimum solution
- Example: Always-on face recognition
- Need for custom offline NN compilers
- Example: Tensilica® Xtensa Neural Network Compiler

## Online

- Varied framework support
- Dynamic application development
- Need most convenient porting solution
- Example: Custom NN app on mobile phone
- Need for dynamic NN compilers
- Example: Android Neural Network API

# Offline Software Solution:
# Tensilica® Xtensa Neural Network Compiler (XNNC)

embedded
VISION
SUMMIT
2018

**Neural Network Descriptor**
(Caffe, TensorFlow, TensorFlow Lite)

Xtensa® Neural Network Compiler

Xtensa C/C++ Compiler

NN Runtime

NN Function Library
(lib "XICNN")

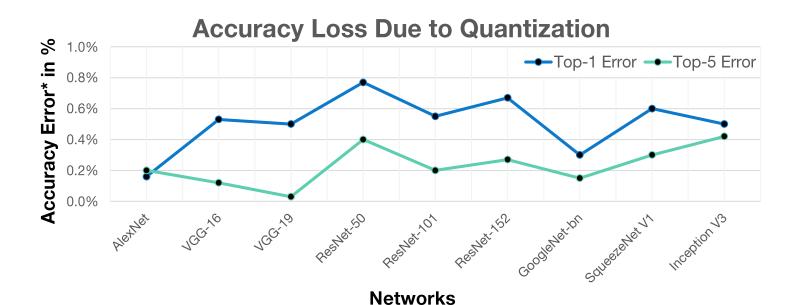DMA Manager (libidma)

Tensilica® Vision P6/Q6/C5 DSP

**Compute & Bandwidth Optimized Code Generator**

- Custom Quantization to 8b data & weights

- Use of target specific optimized NN Library with convolution and non-convolution layers

- Inclusion of performance enhancement features like selection of most optimum library function, kernel fusion, kernel rejection, DMA & tile management

User Code
Tensilica® Compiler / Tool
Tensilica® SW Library / Runtime
Tensilica® Vision DSPs
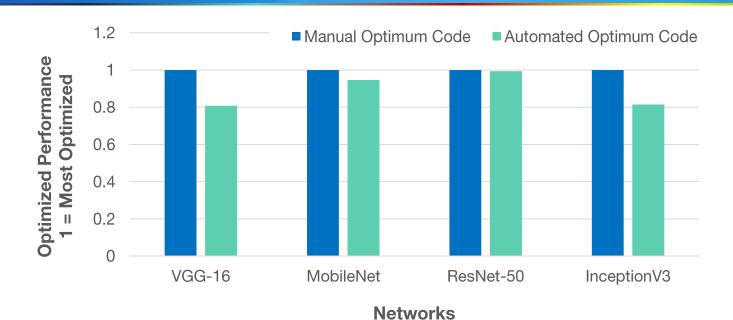
**Accuracy Loss Due to Quantization**

Vision DSP's use of 8b fixed point for AI processing has negligible accuracy impact

*Code Generated through XNNC and tested over 50K Images

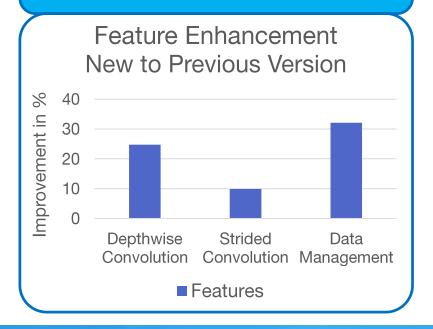# XNNC Highlight: Optimum Code Generation

Automated generated code is close to optimum implementation resulting in speedy time to market

*Code Generated through XNNC for Vision P6 DSP

cādence®

# XNNC Highlights: Continuous Feature Update



embedded
VISION
SUMMIT
2018

## Continuous Investment in Optimization

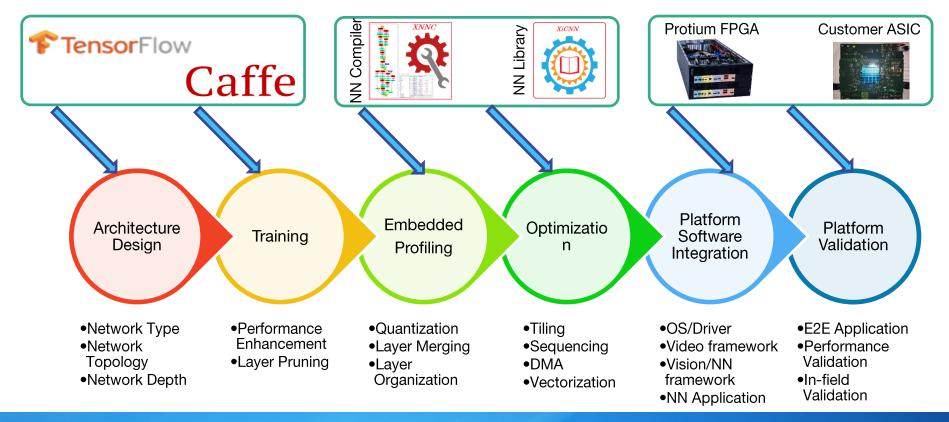**Feature Enhancement New to Previous Version**



## Continuous Addition of New Layers, Framework and Features

- Wide support for framework descriptors: Caffe, TensorFlow, TensorFlow Lite

- Wide support of network types: classification, object detection, segmentation, recurrent, regression

- Custom layer

- Sparsity utilization for bandwidth and compute optimizations

# XNNC with Dynamic Deployment: Time to Market Process

TensorFlow
Caffe

NN Compiler — XNNC
NN Library — XiCNN

Protium FPGA
Customer ASIC

Architecture Design → Training → Embedded Profiling → Optimization → Platform Software Integration → Platform Validation

**Architecture Design**
- Network Type
- Network Topology
- Network Depth

**Training**
- Performance Enhancement
- Layer Pruning

**Embedded Profiling**
- Quantization
- Layer Merging
- Layer Organization

**Optimization**
- Tiling
- Sequencing
- DMA
- Vectorization

**Platform Software Integration**
- OS/Driver
- Video framework
- Vision/NN framework
- NN Application

**Platform Validation**
- E2E Application
- Performance Validation
- In-field Validation

**TensorFlow**
Trained Floating Point Model of MobileNet V1

↓

**XNNC**
Auto Generated Quantized, Optimized Source Code

↓

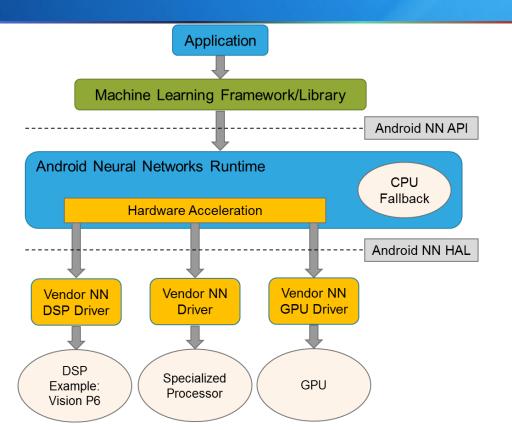Integrated into Vision P6 DSP based on Dreamchip SDK

↓

Verified and validated end-to-end deployment

Several man-months reduced to less than a couple of days

# Online Software Solution: Android Neural Network API



Application

Machine Learning Framework/Library

Android NN API

Android Neural Networks Runtime

CPU Fallback

Hardware Acceleration

Android NN HAL

Vendor NN DSP Driver

Vendor NN Driver

Vendor NN GPU Driver

DSP Example: Vision P6

Specialized Processor

GPU

Courtesy: https://developer.android.com/ndk/guides/neuralnetworks/

## **Easy Deployment of NN on Android Devices**

- Called by ML Libraries or frameworks

- Efficiently distributes workload across available on-device processors

# Android NN Highlights: Dynamic APP Development

## Enabling App Development on Android Devices

- Apps directly use higher level frameworks to deploy trained model on device

- Specialized inferencing engines enables highest performance at lowest power penalty

## Real-Time Optimized Execution on Tensilica® Vision P6 DSP

- Executes graph/sub-graph/layer

- Gets best runtime optimization using tile, DMA management, data rearrangement

- Use of hand-optimized ML library

# Conclusion

- **<u>Trend</u>**: See, hear and speak more clearly with On-Device AI

- **<u>Software</u>**: Need for static and dynamic ecosystem for development and deployment of trained networks

- **<u>Example</u>**:

  - Tensilica® Xtensa Neural Network Compiler

  - Android Neural Network API

See our AI demos at booth 200. Visit our webpage ip.cadence.com/vision.

- Cadence Web Page
  - https://ip.cadence.com

- Android NN
  - https://developer.android.com/ndk/guides/neuralnetworks/index.html