# Building the Adaptable, Intelligent World

# Machine Learning Suite

Kamran Khan

Sr Product Manager, AI and ML

**≡ XILINX.**

Deep Learning explores the study of algorithms that can learn from and make predictions on data

Deep Learning is Re-Defining many Applications

Cloud Acceleration

Security
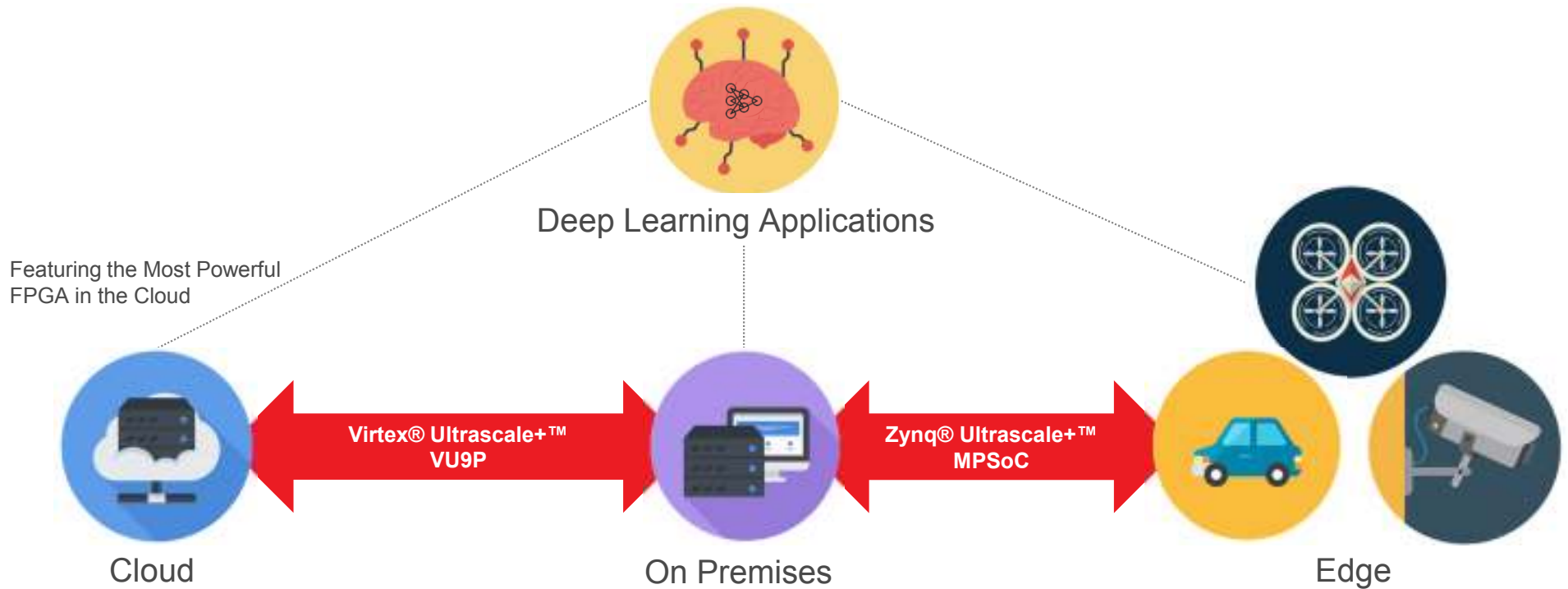
Ecommerce Social

Financial

Surveillance

Industrial IOT

Medical Bioinformatics

Autonomous Vehicles

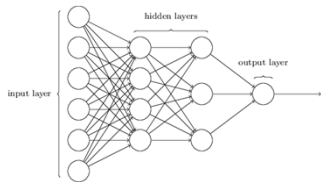# Accelerating AI Inference into Your Cloud Applications



Deep Learning Applications

Featuring the Most Powerful
FPGA in the Cloud

**Virtex® Ultrascale+™
VU9P**

**Zynq® Ultrascale+™
MPSoC**

Cloud

On Premises

Edge

# Overlay Architecture
# Custom Processors Exploiting Xilinx FPGA Flexibility

➤ Customized overlays with ISA architecture for optimized implementation

➤ Easy plug and play with Software Stack



MLP Engine
Scalable sparse and dense implementation



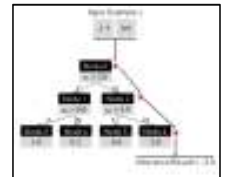xDNN – CNN Engine for Large 16 nm Xilinx Devices
Deephi DPU – Flexible CNN Engine with Embedded Focus
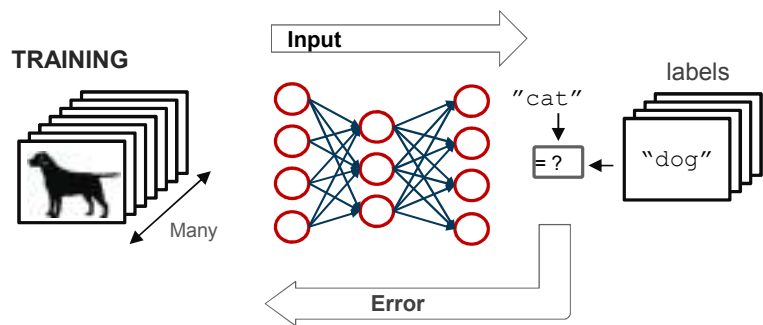CHaiDNN – HLS based open source offering
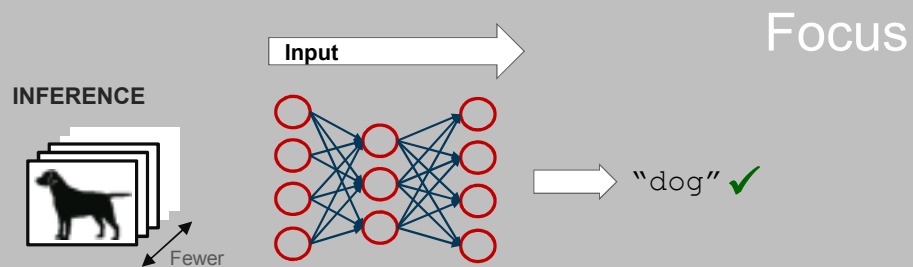


Deephi ESE
LSTM Speech to Text engine



Random Forest
Configurable RF classification

XILINX.

# Machine Learning **Inference** is Xilinx Focus


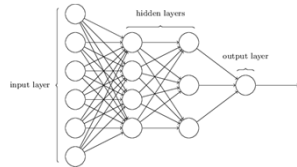
**Training**: Process for machine to "learn" and optimize model from data

**Inference**: Using trained models to predict/estimate outcomes from new observations in efficient deployments

# Deep Learning Models







**Multi-Layer Perceptron**

- Classification
- Universal Function Approximator
- Autoencoder

**Convolutional Neural Network**

- Feature Extraction
- Object Detection
- Image Segmentation

**Recurrent Neural Network**

- Sequence and Temporal Data
- Speech to Text
- Language Translation

**Classification**



"Dog"
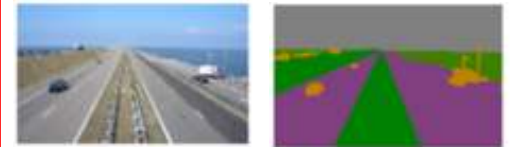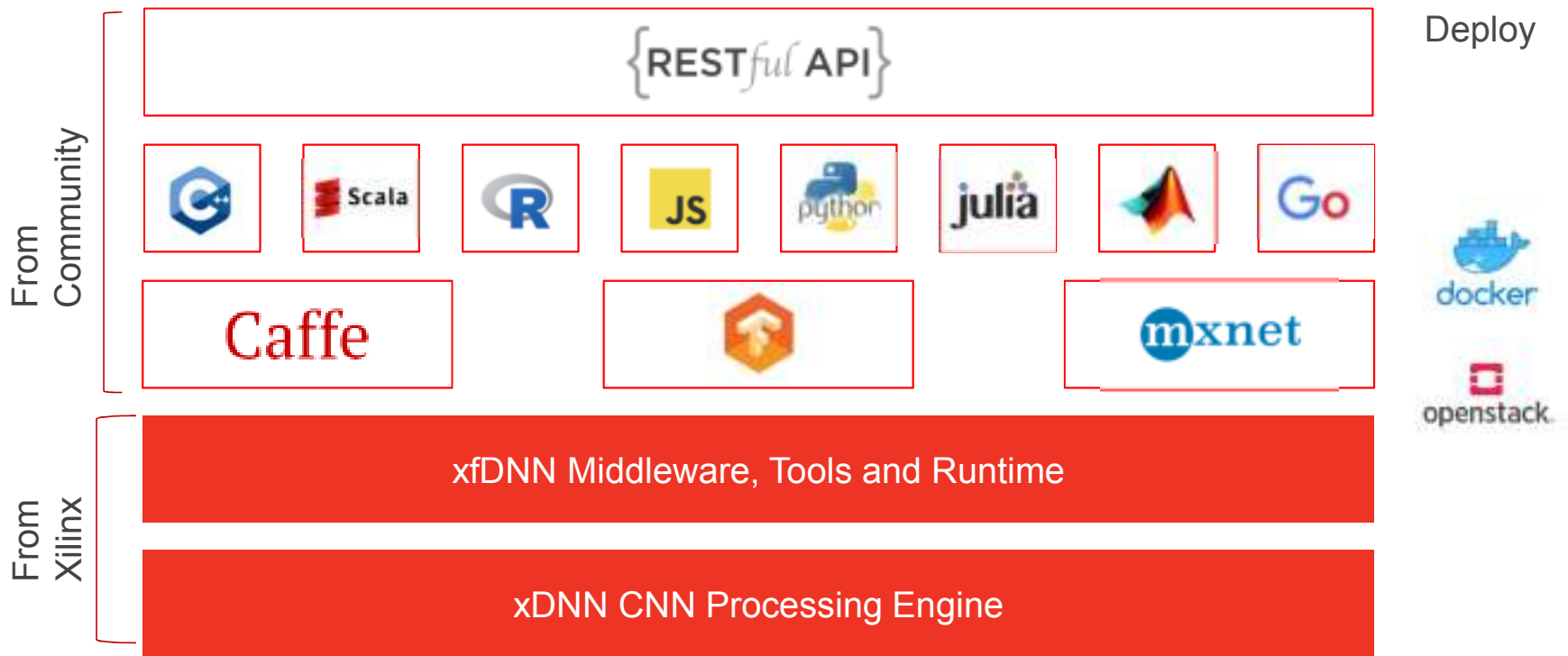
**Object Detection**



**Segmentation**

XILINX

# Seamless Deployment with Open Source Software



Deploy

From Community

{RESTful API}

xfDNN Middleware, Tools and Runtime

xDNN CNN Processing Engine

From Xilinx

*TensorFlow Q4 2017

XILINX

# xfDNN flow

CNTK

Caffe2

PyTorch

Tensorflow

MxNet

Caffe

**FRONTEND**

Framework Tensor Graph to Xilinx Tensor Graph

xfDNN Tensor Graph Optimization

ONNX

xfDNN Compiler

xfDNN Compression

Model Weights

Calibration Set

Image

xfDNN Runtime
*(python API)*

CPU Layers

FPGA Layers

*https://github.com/Xilinx/ml-suite*

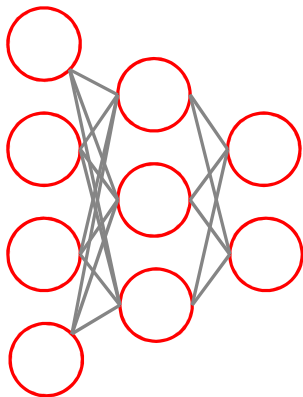**XILINX**

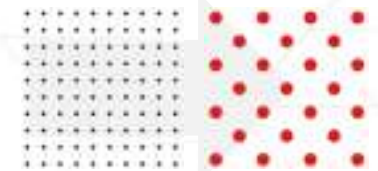# xfDNN Inference Toolbox

## Graph Compiler



- Python tools to quickly compile networks from common Frameworks – Caffe, MxNet and Tensorflow
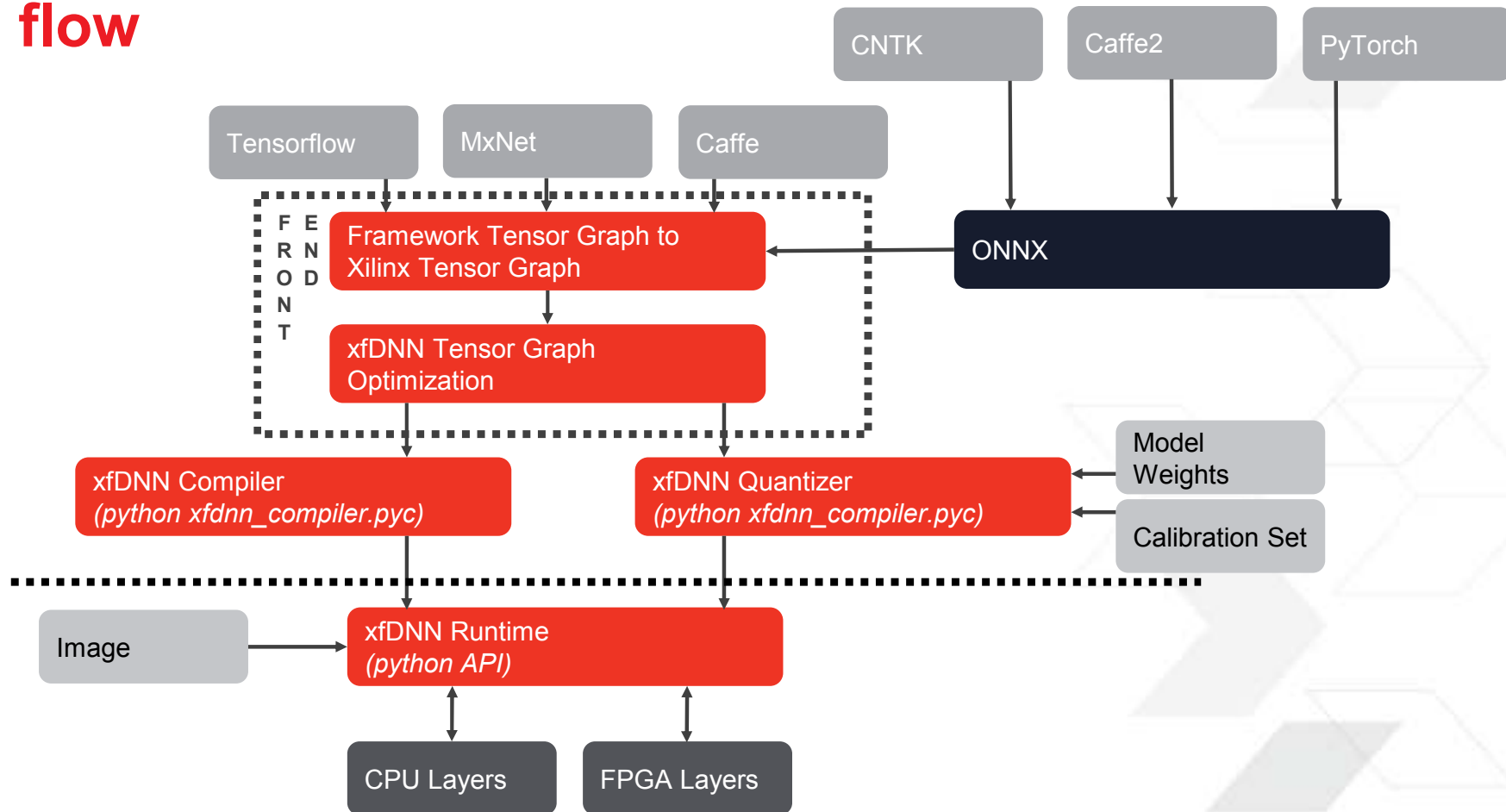
## Network Optimization



- Automatic network optimizations for lower latency by fusing layers and buffering on-chip memory

## xfDNN Quantizer



- Quickly reduce precision of trained models for deployment
- Maintains 32bit accuracy at 8 bit within 2%

XILINX.

# xfDNN flow



CNTK

Caffe2

PyTorch

Tensorflow

MxNet

Caffe

**FRONT END**

Framework Tensor Graph to Xilinx Tensor Graph

xfDNN Tensor Graph Optimization

ONNX

xfDNN Compiler
*(python xfdnn_compiler.pyc)*

xfDNN Quantizer
*(python xfdnn_compiler.pyc)*

Model Weights

Calibration Set

Image

xfDNN Runtime
*(python API)*

CPU Layers

FPGA Layers

*https://github.com/Xilinx/ML-Development-Stack-From-Xilinx*

**XILINX**

# xfDNN Graph Compiler

Pass in a Network

xfDNN
Graph Compiler

Microcode for xDNN is Produced

**XILINX**

# xfDNN Network Deployment



"One Shot" Network Deployment

**Fused Layer Optimizations**
- Compiler can merge nodes
  - (Conv or EltWise)+Relu
  - Conv + Batch Norm
- Compiler can split nodes
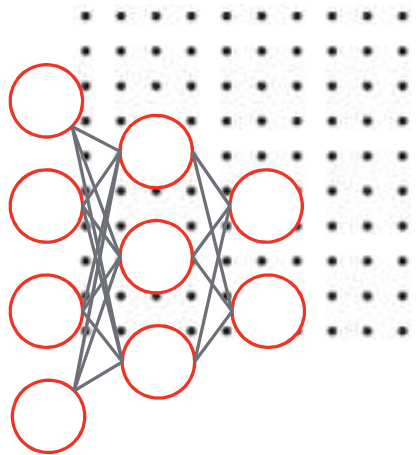  - Conv 1x1 stride 2 -> Maxpool+Conv 1x1 Stride 1

**On-Chip buffering reduces latency and increases throughput**
- xfDNN analyzes network memory needs and optimizes scheduler
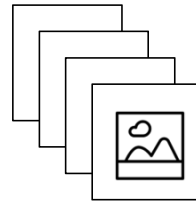  - For Fused and "One Shot" Deployment

**"One Shot" deploys entire network to FPGA**
- Optimized for fast, low latency inference
- Entire network, schedule and weights loaded only once to FPGA

# xfDNN Quantizer: Fast and Easy

1) Provide FP32 network and model
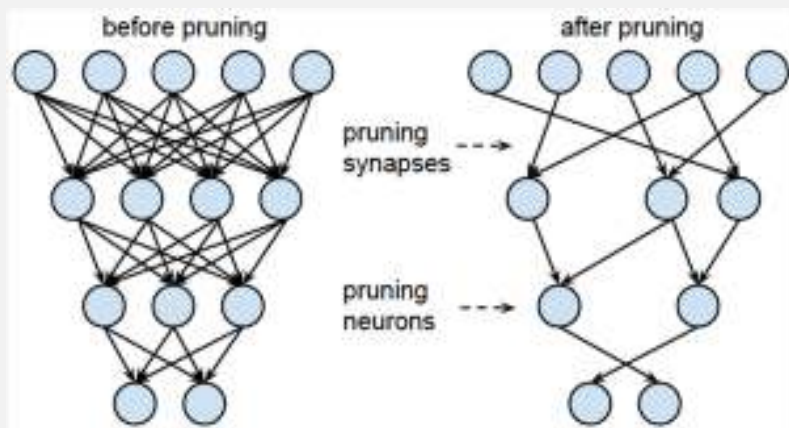   - E.g., prototxt and caffemodel

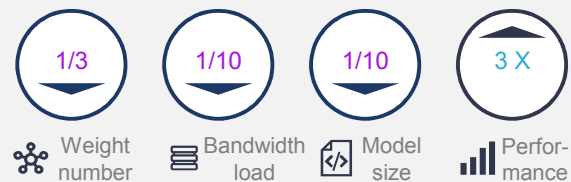2) Provide a small sample set, no labels required
   - 16 to 512 images

Int$_8$

3) Specify desired precision
   - Quantizes to <8 bits to match Xilinx's DSP

**XILINX**

# Xilinx Pruning Overview

**Deep compression**
**Makes algorithm smaller and lighter**


before pruning — after pruning
pruning synapses
pruning neurons

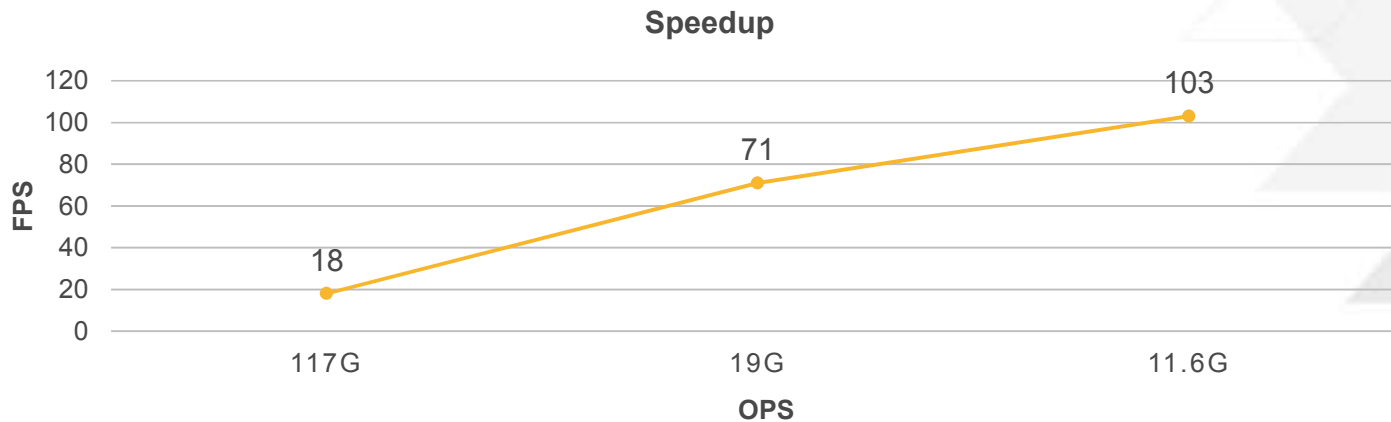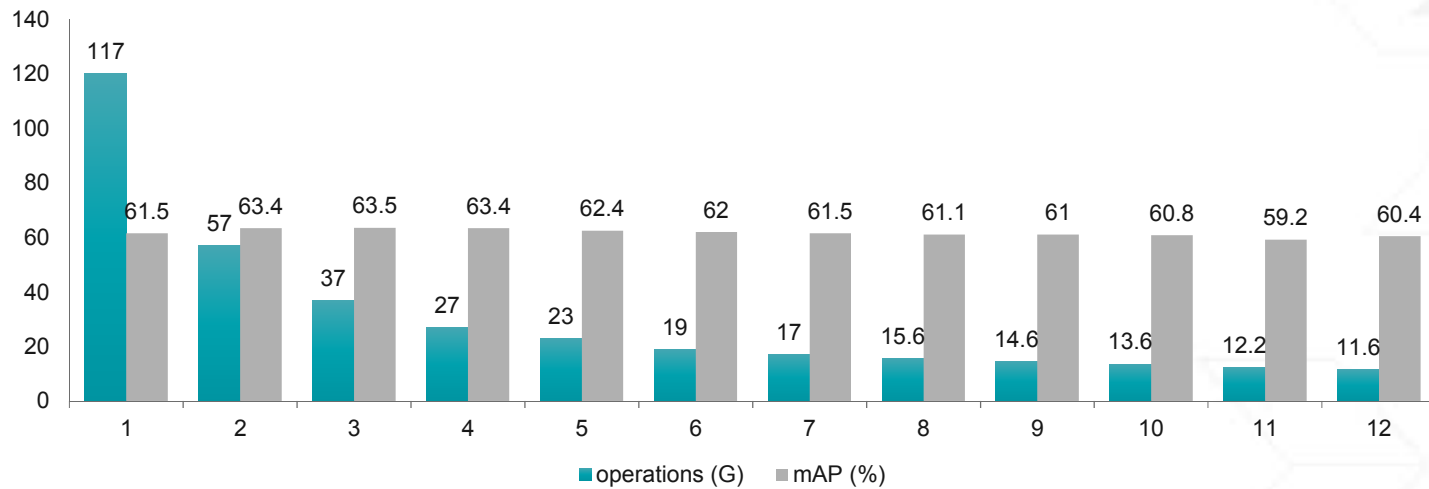| Highlight | 1/3 | 1/10 | 1/10 | 3 X |
|-----------|-----|------|------|-----|
| | Weight number | Bandwidth load | Model size | Performance |

**Compression efficiency**

Deep Compression Tool can achieve significant compression on **CNN** and **RNN**

**Accuracy**

Algorithm can be **compressed 7 times without losing accuracy** under SSD object detection framework

XILINX.

# Pruning Example - SSD



Bar chart — operations (G) and mAP (%):

| # | operations (G) | mAP (%) |
|---|---|---|
| 1 | 117 | 61.5 |
| 2 | 57 | 63.4 |
| 3 | 37 | 63.5 |
| 4 | 27 | 63.4 |
| 5 | 23 | 62.4 |
| 6 | 19 | 62 |
| 7 | 17 | 61.5 |
| 8 | 15.6 | 61.1 |
| 9 | 14.6 | 61 |
| 10 | 13.6 | 60.8 |
| 11 | 12.2 | 59.2 |
| 12 | 11.6 | 60.4 |

■ operations (G)   ■ mAP (%)

## Speedup



| OPS | FPS |
|---|---|
| 117G | 18 |
| 19G | 71 |
| 11.6G | 103 |

**XILINX.**

# Supported DNN (Deep Neural Network) by Applications

| Application | NTT Request | Function | Algorithm |
|---|---|---|---|
| Face | | Face detection | SSD, Densebox |
| | | Landmark Localization | Coordinates Regression |
| | | Face recognition | ResNet + Triplet / A-softmax Loss |
| | | Face attributes recognition | Classification and regression |
| Pedestrian | 1 | Pedestrian Detection (Crowd Volume) | SSD |
| | | Pose Estimation | Coordinates Regression |
| | | Person Re-identification | ResNet + Loss Fusion |
| Video Analytics | 1 | Object detection | SSD, RefineDet |
| | | Pedestrian Attributes Recognition | GoogleNet |
| | | Car Attributes Recognition | GoogleNet |
| | 1 | Car Logo Detection | DenseBox |
| | 1 | Car Logo Recognition | GoogleNet + Loss Fusion |
| | 1 | License Plate Detection | Modified DenseBox |
| | 1 | License Plate Recognition | GoogleNet + Multi-task Learning |
| ADAS/AD | | Object Detection | SSD, YOLOv2, YOLOv3 |
| | | 3D Car Detection | F-PointNet, AVOD-FPN |
| | | Lane Detection | VPGNet |
| | | Traffic Sign Detection | Modified SSD |
| | | Semantic Segmentation | FPN |
| | | Drivable Space Detection | MobilenetV2-FPN |
| | | Multi-task (Detection+Segmentation) | Deephi |

XILINX

# xDNN Process Engine

# Rapid Feature and Performance Improvement

**xDNN-v1**
**Q4CY17**

- Array of Accumulator
- Int16 (Batch=1) and Int8 (Batch=2) support
- Instructions: Convolution, ReLU, Pool, Elementwise
- Flexible kernel size(square) and strides
- 500 MHz

**xDNN-v2**
**Q2CY18**

- All xDNN-v1 Features
- DDR Caching: Larger Image size
- New Instructions: Depth-wise Convolution, De-convolution, Up-sampling
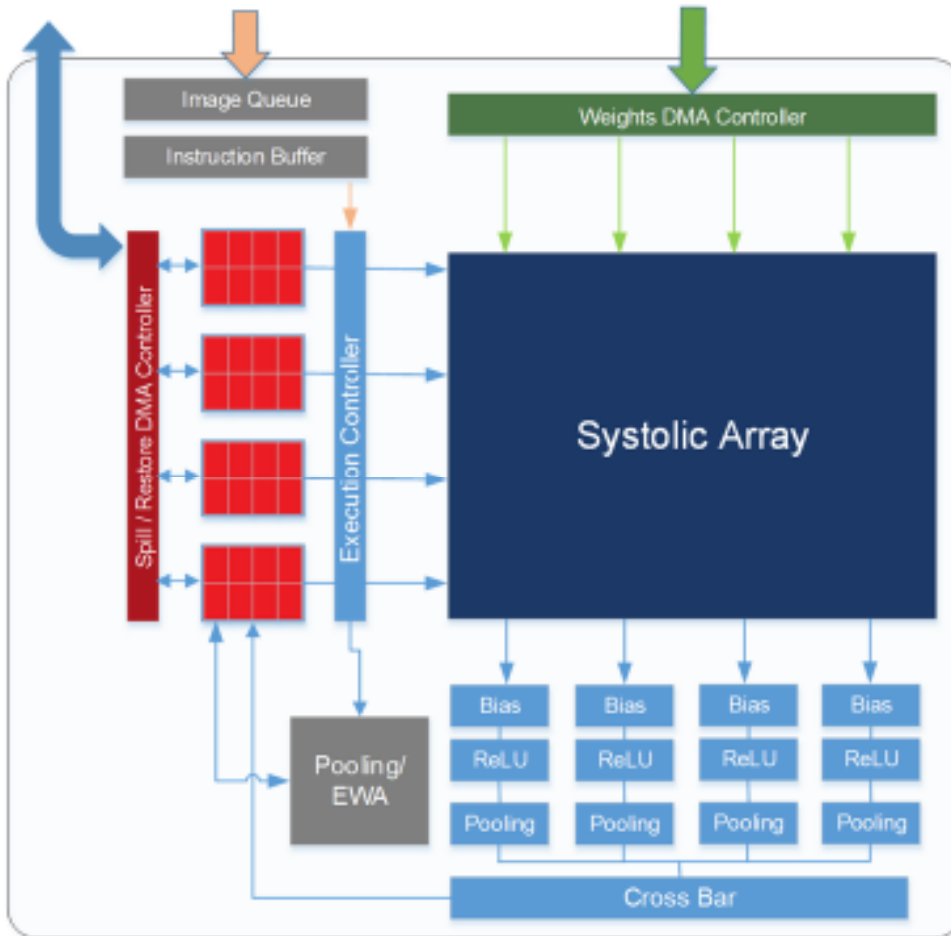- Rectangular Kernels
- 500 MHz

**xDNN-v3**
**Q4CY18**

- New Systolic Array Implementation: 2.2x lower latency
- Instruction Level Parallelism – non-blocking data movement
- Batch=1 for Int8 – lower latency
- Feature compatible with xDNN-v2
- 720+ MHz

XILINX

# XDNN v3 Feature Set

| Features | | | Description | |
|---|---|---|---|---|
| **Supported Operations** | **Convolution / Deconvolution / Convolution Transpose** | Kernel Sizes | W: 1-15; H:1-15 | |
| | | Strides | W: 1,2,4,8; H: 1,2,4,8 | |
| | | Padding | Same, Valid | |
| | | Dilation | Factor: 1,2,4 | |
| | | Activation | ReLU/pReLU | |
| | | Bias | Value Per Channel | |
| | | Scaling | Scale & Shift Value Per Channel | |
| | **Max Pooling** | Kernel Sizes | W: 1-15; H:1-15 | |
| | | Strides | W: 1,2,4,8; H: 1,2,4,8 | |
| | | Padding | Same, Valid | |
| | **Avg Pooling** | Kernel Sizes | W: 1-15; H:1-15 | |
| | | Strides | W: 1,2,4,8; H: 1,2,4,8 | |
| | | Padding | Same, Valid | |
| | **Element-wise Add** | Width & Height must match; Depth can mismatch. | | |
| | **Memory Support** | On-Chip Buffering, DDR Caching | | |
| | **Expanded set of image sizes** | Square, Rectangular | | |
| | **Upsampling** | Strides | Factor: 2,4,8,16 | |
| **Miscellaneous** | **Precision** | Int16-bit or Int8-bit | | |

**XILINX**

# Xilinx DNN Processor (xDNN)



> Configurable Overlay Processor

> DNN Specific Instruction Set
>> Convolution, Max Pool etc.

> Any Network, Any Image Size

> High Frequency & High Compute Efficiency

> Compile and run new networks

© Copyright 2018 Xilinx

XILINX.

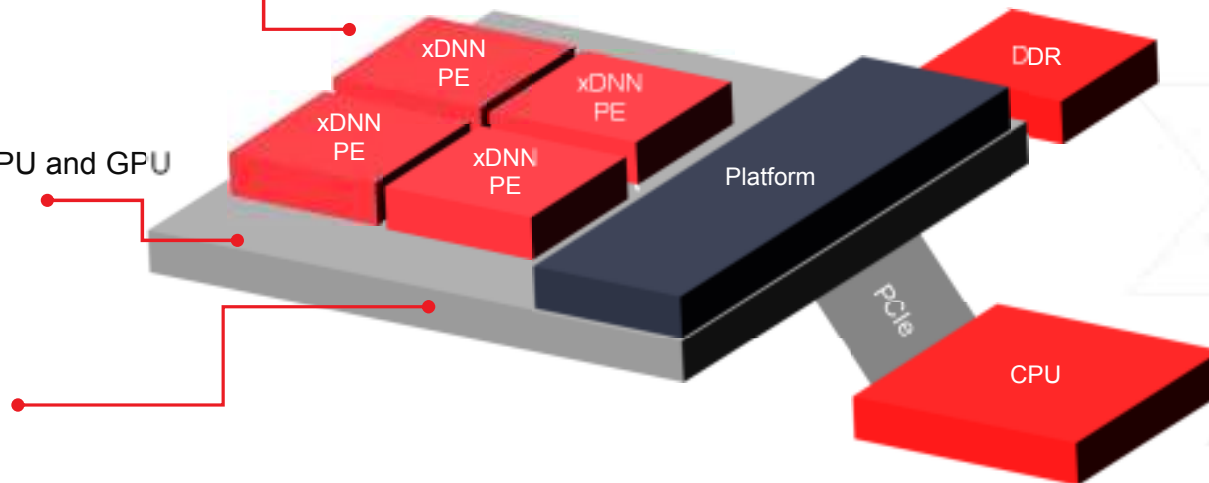# ML Suite Overlays with xDNN Processing Engines

**Adaptable**

> AI algorithms are changing rapidly
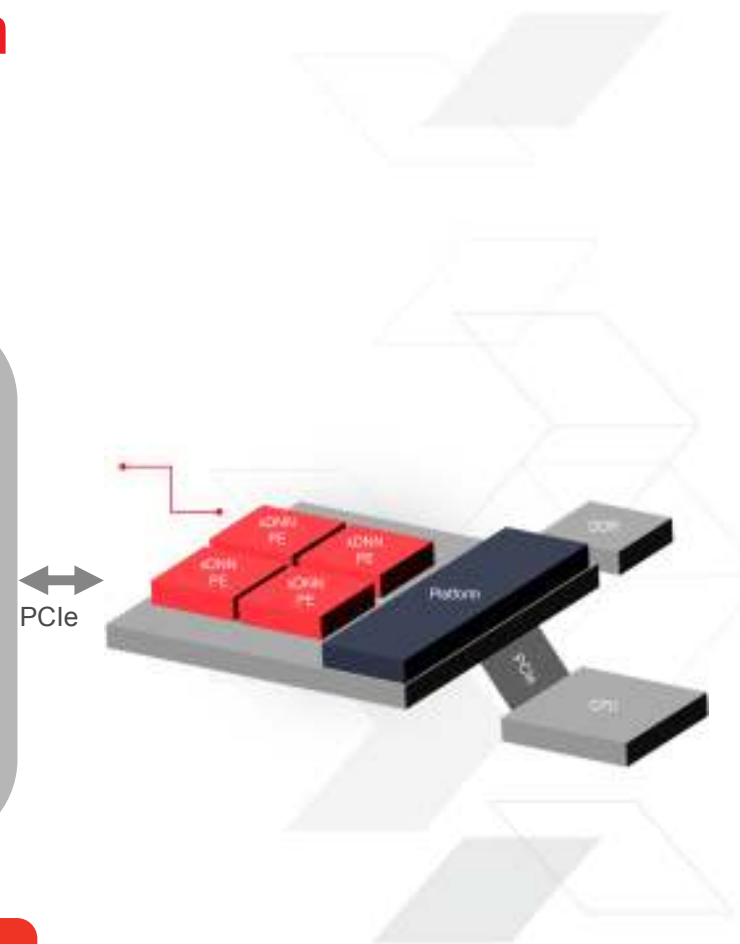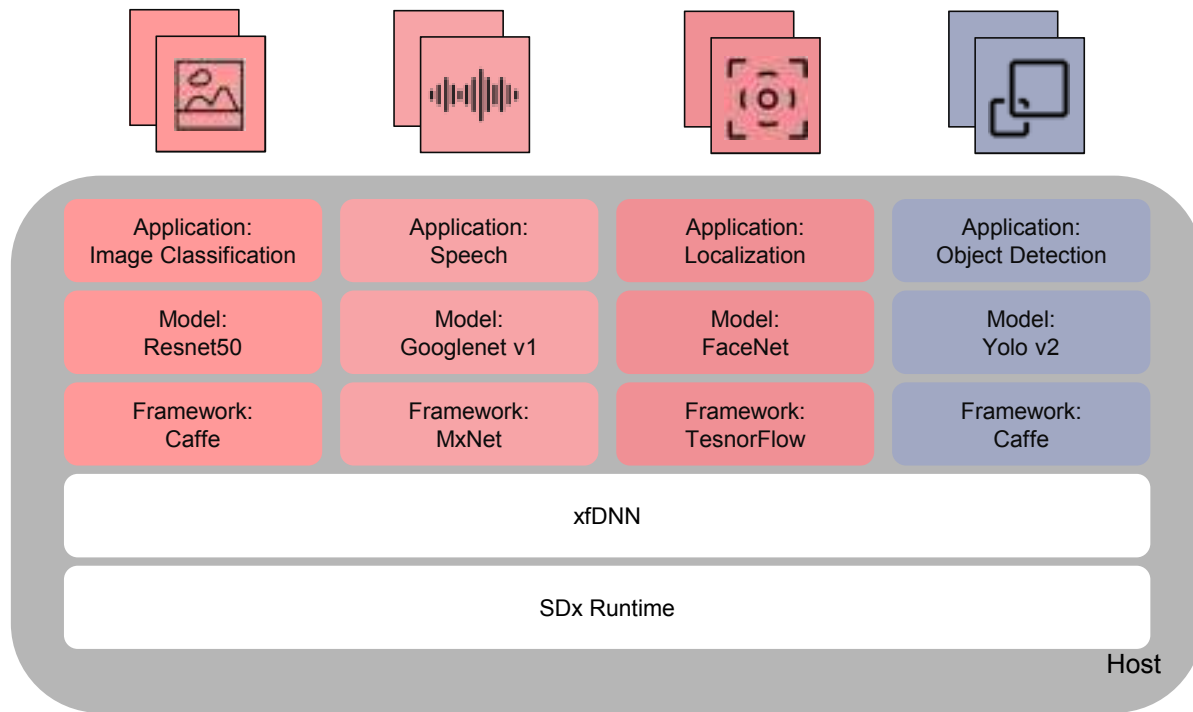
> Adjacent acceleration opportunities

**Realtime**

> 10x Low latency than CPU and GPU

> Data flow processing

**Efficient**

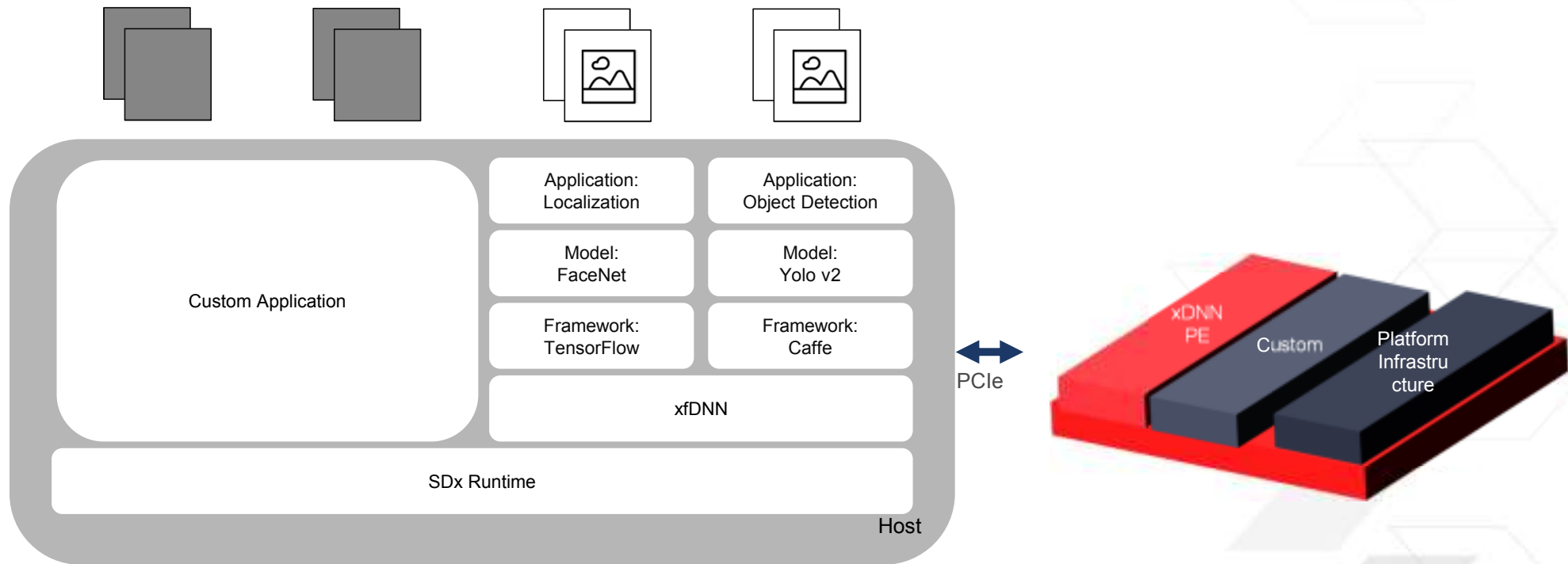> Performance/watt

> Low Power

# Flexible: Multi-Network Configuration



| Application:<br>Image Classification | Application:<br>Speech | Application:<br>Localization | Application:<br>Object Detection |
|---|---|---|---|
| Model:<br>Resnet50 | Model:<br>Googlenet v1 | Model:<br>FaceNet | Model:<br>Yolo v2 |
| Framework:<br>Caffe | Framework:<br>MxNet | Framework:<br>TesnorFlow | Framework:<br>Caffe |

xfDNN

SDx Runtime

Host

PCIe

**1 FPGA Provides 4 Virtual Accelerators
For Real Time Deep Learning**

**XILINX.**

# Flexible: Bring Your own IP!



Custom Application

| Application: Localization | Application: Object Detection |
|---|---|
| Model: FaceNet | Model: Yolo v2 |
| Framework: TensorFlow | Framework: Caffe |
| xfDNN | |

SDx Runtime

Host

PCIe

xDNN PE

Custom

Platform Infrastructure

## Integrate Custom Applications Directly with xDNN Processing Engines

**XILINX**

# Alveo – Breathe New Life into Your Data Center

**16nm UltraScale™ Architecture**

**NIMBIX** Cloud Deployed

**Off-Chip Memory Support**
- Max Capacity: 64GB
- Max Bandwidth: 77GB/s

**Cloud ↔ On-Premise Mobility**

**Internal SRAM**
- Max Capacity: 54MB
- Max Bandwidth: 38TB/s

**Ecosystem of Applications**
- Many available today
- More on the way

**PCIe Gen3x16 PCIe Gen4x8 w/ CCIX**

**Server OEM Support**
- Major OEMs in Qualification

**HBM2 Memory Support**
- Max Capacity: 8GB
- Max Bandwidth: 460GB/s

**SDAccel** Environment
**Accelerate Any Application**
- IDE for compiling, debugging, profiling
- Supports C/C++, RTL, and OpenCL

XILINX

**U200**

892K
LUTs

35MB
Internal SRAM
Capacity

31TB/s
Internal SRAM
Bandwidth

3100$_{img/s}$
CNN Throughput*

**U250**

1,341K
LUTs

54MB
Internal SRAM
Capacity

38TB/s
Internal SRAM
Bandwidth

4100$_{img/s}$
CNN Throughput*

*Low-latency GoogLeNet v1

XILINX

# It's All About the Applications



**Database**
- BLACKLYNX — Elasticsearch
- ALGO-LOGIC — Low Latency KVS
- bigstream — ETL, Streaming, SQL Analytics
- VITESSE DATA — Deepgreen DB
- XELERA — Spark-MLlib
- Titan — Hyperion 10G RegEx

**Machine Learning**
- XILINX — ML Suite for Inference
- Mipsology — Image Classification
- SumUp — Text Analysis

**Video**
- NGCODEC — ABR Transcoding
- skreens — Experience Engine
- CTACCEL — Image Processing

**Financial**
- MAXELER Technologies — Real-Time Risk Dashboard
- MAXELER Technologies — SIMM Calculation

**HPE & Life Sciences**
- Falcon COMPUTING — Accelerated Genomics Pipelines

## Application Ecosystem Continues to Grow

*Alveo Applications Directory*

>> 27

XILINX

# Solution Stack



**Accelerated Solutions**

**DEVELOPERS**

**100%**
Growth of Published Applications

**Hundreds**
of Developers Trained

RTL, C, C++, OpenCL

**End user**

Xilinx ML Suite

**Framework, API, Python/Java/C++ Programmability**

| HPC &LIFE SCIENCES | FINANCIAL | VIDEO |
|---|---|---|
| MACHINE LEARNING | | DATABASE |

Solutions Xilinx ISVs

**Developer Package**

SDAccel Environment

**Platforms**

FPGA as a Service (FaaS)

Cloud

On-premise

Platform

**XILINX.**

# Orchestration (Kubernetes, OpenShift, Etc)

**Migrate the components to the right hardware resource and manage the container communication across the cluster.**

© Copyright 2018 Xilinx

XILINX

# ML Suite Performance Roadmap

**ACAP**

7x Performance Improvement

**Model Compression, Low and Mixed precision**

Img/Sec
GoogleNet v1, Batch < 2

- 9000 — V100 Batch = 128
- 8000
- 7000
- 6000
- 5000
- 4000 — VCU1526 (VU13P)
- 3000 — VCU1525 (VU9P) / VCU155x (VU37P)
- 2000
- 1000 — VCU1525 (VU9P) / P4 / V100
- K80 B8 / Xeon

Today | Mar'18 | Jun'18 | Sep'18 | Dec'18 | Mar'19 | Jun'19

**Q4FY18** | **Q1FY19** | **Q2FY19** | **Q3FY19** | **Q4FY19** | **Q1FY20** | **Q2FY20**

**Legend:**
- xDNNv2, INT8
- xDNNv3, INT8
- Versal AI Core, INT8

**XILINX**

# xDNN YOLO v2 Performance – Image Size 608x608

© Copyright 2018 Xilinx

XILINX

# Key Differentiating Value vs. Nvidia

Xilinx Enables . . . .

> **Highest Throughput <u>WITH</u> Low Latency (i.e. low batch)**
>> <u>AND</u> best Perf / Watt

> **Accelerate the <u>Whole Application</u> for Broad Range of Workloads**

> **Innovate Faster than the Silicon Cycle with Adaptable Hardware**

Real-Time Inference: GoogleNetv1 (Batch=1)



*Execution Time*

CPU + GPU | Pre | ML | Post

CPU + FPGA | Pre | ML | Post

Workload Evolution

GPGPU or Fixed Function Accelerator Cards — New Hardware Design — New Hardware Design

Alveo Adaptable Accelerator Cards — Same hardware: Reconfigurable for workload optimization

Future Proof

XILINX

# What About Batching?

## Fundamental to GPU Architecture

(Software Defined Data Flow)

**Batching**: Loading up lots of similar Data Sets
- Keep CUDA cores busy
- Hide some memory latency
- Create better SIMT efficiency

| Input 1 | GPU |
|---------|-----|
| Input 2 | Result 1 |
| Input 3 | Result 2 |
| Input 4 | Result 3 |
|         | Result 4 |

High Throughput OR Low Latency

## Not Required for FPGA / ACAP

(Hardware Defined Data Flow)

Independent of Data Set count
- Custom HW kernels
- Custom Memory Hierarchy
- HW pipeline data flow

| Input 1 | FPGA/ACAP |
|---------|-----------|
| Input 2 | Result 1 |
| Input 3 | Result 2 |
| Input 4 | Result 3 |
|         | Result 4 |

High Throughput AND Low Latency

**XILINX.**

# A Batching Example: Image Classification

GPU

DRAM

L2

"Dog"

Time

Memory  Compute

Batch = 1

Low Throughput
Low Latency

XILINX.

# A Batching Example: Image Classification

GPU



L2

"Dog" "Cat" "Truck" "Duck" "Cat" "Dog" "Duck" "Truck"

FPGA / ACAP

xDNN Kernels

DDR

URAM BRAM

"Cat"
"Duck"
"Truck"
"Dog"

Time

Memory   Compute

Batch = 1

Low Throughput
Low Latency

Memory   Compute

Batch = 64

High Throughput
High Latency

High Throughput

AND

Low Latency

XILINX.

# Overview of Compressed LSTM



Our algorithm, software and hardware co-design flow for RNN(LSTM) acceleration

# LSTM Solution Features

## Features

Innovative full-stack acceleration solution for deep learning in acoustic speech recognition (ESE: best paper of FPGA2017)

- Support both unidirectional and bi-directional LSTM acceleration on FPGA for model inference
- Support CNN layers, Fully-Connected (FC) layers, Batch Normalization layers and varieties of activation functions such as Sigmoid, Tanh and HardTanh
- Support testing for both performance comparison of CPU/FPGA and single sentence recognition
- Supporting user's own test audio recognition (English, 16kHz sample rate, no longer than 3 seconds)

| | |
|---|---|
| Usage | Hardware PCIE interface, software API |
| Supported layer | CNN、 uni/bi-directional LSTM(P)、 FC、 BN |
| LSTM layer number | According to the requirements and source |
| Channel number | According to the requirements and source |
| Quantization | 16bit |
| Maximum input of LSTM | 1024 |
| Maximum size of LSTM | 2048 |
| Density of LSTM | Any, typically 10%~20% |
| Peephole in LSTM | Selectable |
| Projection in LSTM | Selectable |
| Activation function | Sigmoid、 Tanh、 HardTanh |

Note: It is available to change the hardware configuration for different requirements such as layer number and model size

**Result**

ASR system

- Decoder
- model
- Softmax — 1 layer only used in training
- Fully connected — 1 layer
- LSTM — 5 layers uni- or bi- directional input = 672, cell = 800
- CNN — 2 layers with BN and Hardtanh
- Feature extraction

**Audio**

XILINX

# Competitive Performance



|  | Sogou Case | DeepSpeech2 |
|---|---|---|
| Model | 3-layer LSTM | 2-layer CNN + 5-layer bidirectional LSTM |
| LSTM Density | 17.76% | 15% |
| (Bi)LSTM Dimension | Input 153, cell 2048, output 1024 | Input 672, cell 800, output 800 |
| CPU Performance | 18155.22ms (E5-2690 V4) | 118.31ms (E5-2686 V4) [BLSTM Only] |
| FPGA Implementation | 32 channels on ku115 @ 250MHz | 1 channel on vcu1525 @ 310MHz (u200@330mhz) |
| Application | Part of ASR in Sogou | End2End ASR |

**XILINX.**

# Integrated Edge – Cloud Roadmap

|  | Edge/Embedded | Cloud/DC |
|---|---|---|

**Models**
- 20+ pruned / customized / basic models
- Deephi Pruning

Coming to ML Suite at XDF

**Software Stack**

| Edge/Embedded | | Cloud/DC |
|---|---|---|
| DECENT | | DECENT / xfDNN Quantizer |
| DNNDK | SDSoC \| SDAccel | xfDNN Compiler |
| DNNDK | | xfDNN Runtime |

**FPGA IP**

| DPU | LSTM | xDNN |
|---|---|---|

**Platforms**

Z7020 Board   Z7020 SOM   ZU2/3 SOM   ZU2/3 Card

Xilinx U200, U250, U280

NIMBIX

Alibaba Cloud aliyun.com

HUAWEI

ZU9 Card   ZCU102   ZCU104   Ultra96

XILINX

# Try the Xilinx ML Suite Today

https://github.com/Xilinx/ml-suite

AWS EC2 Marketplace



Nimbix NX5

# Engage with the Community

© Copyright 2018 Xilinx

# Building the Adaptable, Intelligent World

XILINX

# Building the Adaptable, Intelligent World

XILINX.

# A Broader View of ML Benchmarking  (Int8)

## Throughput vs. GPU Batch

### GoogleNetv1 (Int8)



Latency Critical

Latency Insensitive

Images / Sec

9000
8000
7000
6000
5000
4000
3000
2000
1000
0

U250 : < 2ms Fixed

U200 : < 2ms Fixed

P4    P40    T4 (est)          P4    P40    T4 (est)          P4    P40    T4 (est)          P4    P40    T4 (est)

Batch = 1               Sub 2ms                 Sub 7ms                High Batch

## Alveo Delivers Low, Fixed Latency (< 2ms) in ALL Scenarios

>> 44

XILINX.

# Next Step:  Deephi Pruning Techniques

## Throughput vs. GPU Batch

GoogleNetv1 (Int8)



Deephi Proprietary Pruning Increases Performance **30% or More**

**XILINX.**

# Compute Capability

> **INT8 is becoming standard for inference devices**

Peak Performance at Different Compute Data Types



Today

Legend: Xilinx Alveo U200 | Xilinx Alveo U250 | Nvidia Tesla T4 | Huawei Ascend 310 DC Server Card (4x310)

X-axis: MAX 8b/8b, MAX 4b/8b, MAX 4b/4b, MAX 2b/8b, MAX 2b/2b, MAX T/8b, MAX B/8b, MAX B/4b, MAX B/B, XFP7 (1,3,3), XFP8 (1,3,4)

Y-axis: 1, 10, 100, 1,000, 10,000

XILINX

# Compute Capability

> **Trends are moving toward lower precision, e.g. INT4 is emerging**



Peak Performance at Different Compute Data Types

# Compute Capability

> **Mixed weight and activation precisions provide additional optimization for inference while keeping high accuracy**



Peak Performance at Different Compute Data Types

# Compute Capability

> **The future we see: each model will have its own mix of datatype for optimal efficiency and accuracy**



Peak Performance at Different Compute Data Types

# Compute Capability

> **Alveo: Most future proof architecture**

> **Scalable performance for any data type**



Peak Performance at Different Compute Data Types

# Break Through on Peak Performance

- **GPU: Introduce new architectures and silicon**
- **Xilinx: Adapt the break through of emerging domain knowledge**

### GPU Deep Learning Peak Power Efficiency

Mixed FP16 for TensorCore

Native INT8 operation

Legend: Kepler — Maxwell — Pascal — Volta

### FPGA Deep Learning Peak Power Efficiency

ACAP Architectures

Adaptable Precision

INT8 Optimization

Legend: Ultrascale — Ultrascale+ — ACAP

XILINX.

# Infuse Machine Learning with other accelerations



Database

Machine Learning

HPC & Life Sciences

Video

Financial

XILINX

# Custom Deep Learning Pipeline



**Integrate Custom Applications with xDNN. Lower end-to-end latency**

© Copyright 2018 Xilinx

XILINX

# X + ML Focus Applications Summary

Smart City / Cloud Surveillance
• 10x Lower Latency

High Resolution Imaging
• 1.6x Faster

Security / Anomaly / Malware
• 5x Faster

**XILINX.**

# ML Suite Solution Stack

Deep Learning explores the study of algorithms that can learn from and make predictions on data

Deep Learning is Re-Defining many Applications

Cloud Acceleration

Security

Ecommerce Social

Financial

Surveillance

Industrial IOT

Medical Bioinformatics

Autonomous Vehicles

# Accelerating AI Inference into Your Cloud Applications



Deep Learning Applications

Featuring the Most Powerful
FPGA in the Cloud

**Virtex® Ultrascale+™
VU9P**

**Zynq® Ultrascale+™
MPSoC**

Cloud

On Premises

Edge

# Overlay Architecture
# Custom Processors Exploiting Xilinx FPGA Flexibility

➤ Customized overlays with ISA architecture for optimized implementation

➤ Easy plug and play with Software Stack



MLP Engine
Scalable sparse and dense implementation

xDNN – CNN Engine for Large 16 nm Xilinx Devices
Deephi DPU – Flexible CNN Engine with Embedded Focus
CHaiDNN – HLS based open source offering

Deephi ESE
LSTM Speech to Text engine

Random Forest
Configurable RF classification

XILINX

# Deep Learning Models



## Multi-Layer Perceptron

- Classification
- Universal Function Approximator
- Autoencoder

## Convolutional Neural Network

- Feature Extraction
- Object Detection
- Image Segmentation

## Recurrent Neural Network

- Sequence and Temporal Data
- Speech to Text
- Language Translation

**Classification**



"Dog"

**Object Detection**



**Segmentation**

XILINX

# Seamless Deployment with Open Source Software



Deploy

From Community

From Xilinx

RESTful API

Caffe

xfDNN Middleware, Tools and Runtime

xDNN CNN Processing Engine

docker

openstack

*TensorFlow Q4 2017

XILINX

# xDNN Process Engine

# Rapid Feature and Performance Improvement

**xDNN-v1**
**Q4CY17**

- Array of Accumulator
- Int16 (Batch=1) and Int8 (Batch=2) support
- Instructions: Convolution, ReLU, Pool, Elementwise
- Flexible kernel size(square) and strides
- 500 MHz

**xDNN-v2**
**Q2CY18**

- All xDNN-v1 Features
- DDR Caching: Larger Image size
- New Instructions: Depth-wise Convolution, De-convolution, Up-sampling
- Rectangular Kernels
- 500 MHz

**xDNN-v3**
**Q4CY18**

- New Systolic Array Implementation: 2.2x lower latency
- Instruction Level Parallelism – non-blocking data movement
- Batch=1 for Int8 – lower latency
- Feature compatible with xDNN-v2
- 720+ MHz

**XILINX**

# XDNN v3 Feature Set

| Features | | Description | |
|---|---|---|---|
| **Supported Operations** | **Convolution / Deconvolution / Convolution Transpose** | Kernel Sizes | W: 1-15; H:1-15 |
| | | Strides | W: 1,2,4,8; H: 1,2,4,8 |
| | | Padding | Same, Valid |
| | | Dilation | Factor: 1,2,4 |
| | | Activation | ReLU/pReLU |
| | | Bias | Value Per Channel |
| | | Scaling | Scale & Shift Value Per Channel |
| | **Max Pooling** | Kernel Sizes | W: 1-15; H:1-15 |
| | | Strides | W: 1,2,4,8; H: 1,2,4,8 |
| | | Padding | Same, Valid |
| | **Avg Pooling** | Kernel Sizes | W: 1-15; H:1-15 |
| | | Strides | W: 1,2,4,8; H: 1,2,4,8 |
| | | Padding | Same, Valid |
| | **Element-wise Add** | Width & Height must match; Depth can mismatch. | |
| | **Memory Support** | On-Chip Buffering, DDR Caching | |
| | **Expanded set of image sizes** | Square, Rectangular | |
| | **Upsampling** | Strides | Factor: 2,4,8,16 |
| **Miscellaneous** | **Precision** | Int16-bit or Int8-bit | |

**XILINX**

# Xilinx DNN Processor (xDNN)



> Configurable Overlay Processor

> DNN Specific Instruction Set
>> Convolution, Max Pool etc.

> Any Network, Any Image Size

> High Frequency & High Compute Efficiency

> Compile and run new networks

**XILINX**

# xDNN Compiler + Runtime



*https://github.com/Xilinx/ml-suite*

# xDNN v3 Implementation on Alveo U200

> 3 Large 96x16 PEs– 1 in each SLR – 5.2 ML Shell

> Kernels @ 720 MHz/360MHz

| Resource | Count | Utilization |
|----------|-------|-------------|
| LUTs | 658k | 52% |
| DSPs | 5661 | 80% |
| BRAM | 1258 | 58% |
| URAM | 864 | 92% |

**XILINX**

# xDNN v3 Implementation on Alveo U250

> 4 Large 96x16 PEs– 1 in each SLR – standard 5.2 Shell

> Kernels at 700 MHz/350 MHz

| Resource | Count | Utilization |
|----------|-------|-------------|
| LUTs | 876k | 51% |
| DSPs | 7548 | 62% |
| BRAM | 1632 | 61% |
| URAM | 1152 | 90% |

**XILINX.**

# ML Suite Overlays with xDNN Processing Engines

**Adaptable**

> AI algorithms are changing rapidly

> Adjacent acceleration opportunities

**Realtime**

> 10x Low latency than CPU and GPU

> Data flow processing

**Efficient**

> Performance/watt

> Low Power

**XILINX**

# Flexible: Multi-Network Configuration



| Application:<br>Image Classification | Application:<br>Speech | Application:<br>Localization | Application:<br>Object Detection |
|---|---|---|---|
| Model:<br>Resnet50 | Model:<br>Googlenet v1 | Model:<br>FaceNet | Model:<br>Yolo v2 |
| Framework:<br>Caffe | Framework:<br>MxNet | Framework:<br>TesnorFlow | Framework:<br>Caffe |

xfDNN

SDx Runtime

Host

PCIe

**1 FPGA Provides 4 Virtual Accelerators
For Real Time Deep Learning**

**XILINX**

# Flexible: Bring Your own IP!



Custom Application

| Application: Localization | Application: Object Detection |
|---|---|
| Model: FaceNet | Model: Yolo v2 |
| Framework: TensorFlow | Framework: Caffe |

xfDNN

SDx Runtime

Host

PCIe

xDNN PE    Custom    Platform Infrastructure

Integrate Custom Applications Directly
with xDNN Processing Engines

**XILINX**

# xDNN GoogLeNet v1 Performance – Image Size 224x224

**XILINX**

# xDNN YOLO v2 Performance – Image Size 608x608

© Copyright 2018 Xilinx

**XILINX**

# xDNN Real Time Performance

xDNN v3 Throughput  Performance

Legend: ■ U200 xDNNv3   ▨ U200 xDNNv3.1   ■ U250 xDNNv3   ▨ U250 xDNNv3.1

Categories: Googlenetv1 - Int8 - 224, ResNet50 - Int8 - 224, Yolov2 - Int8 - 224

XILINX

# Xilinx Inference Middleware - xfDNN

**XILINX.**

# Unified Simple User Experience from Cloud to XBB

Develop

Publish

Deploy



User Guides    Tutorials    Examples

xfDNN    Apps    xDNN Binaries

User Guides    Tutorials    Examples

xfDNN    Apps    xDNN Binaries

AWS

Launch Instance

User Choice

Download

XILINX

# Xilinx ML Suite

> **ML Suite**
>> Supported Frameworks:
- Caffe
- MxNet
- Tensorflow
- Python Support
- Darknet

>> Jupyter Notebooks available:
- Image Classification with Caffe
- Using the xfDNN Compiler w/ a Caffe Model
- Using the xfDNN Quantizer w/ a Caffe Model

>> Pre-trained Models
- Caffe 8/16-bit
    - GoogLeNet v1
    - ResNet50
    - Flowers102
    - Places365
- Python 8/16-bit
    - Yolov2
- MxNet 8/16-bit
    - GoogLeNet v1

>> xfDNN Tools
- Compiler
- Quantizer



https://github.com/Xilinx/ml-suite

# Seamless Deployment with Open Source Software

# xfDNN flow

CNTK

Caffe2

PyTorch

Tensorflow

MxNet

Caffe

**FRONTEND**

Framework Tensor Graph to Xilinx Tensor Graph

ONNX

xfDNN Tensor Graph Optimization

xfDNN Compiler

xfDNN Compression

Model Weights

Calibration Set

Image

xfDNN Runtime
*(python API)*

CPU Layers

FPGA Layers

*https://github.com/Xilinx/ml-suite*

**XILINX**

# xfDNN Inference Toolbox

## Graph Compiler



- Python tools to quickly compile networks from common Frameworks – Caffe, MxNet and Tensorflow

## Network Optimization



"One Shot" Network Deployment

- Automatic network optimizations for lower latency by fusing layers and buffering on-chip memory

## xfDNN Quantizer



- Quickly reduce precision of trained models for deployment
- Maintains 32bit accuracy at 8 bit within 2%

**XILINX.**

# xfDNN flow



CNTK

Caffe2

PyTorch

Tensorflow

MxNet

Caffe

**FRONT END**

Framework Tensor Graph to Xilinx Tensor Graph

xfDNN Tensor Graph Optimization

ONNX

Model Weights

Calibration Set

xfDNN Compiler
*(python xfdnn_compiler.pyc)*

xfDNN Quantizer
*(python xfdnn_compiler.pyc)*

Image

xfDNN Runtime
*(python API)*

CPU Layers

FPGA Layers

*https://github.com/Xilinx/ML-Development-Stack-From-Xilinx*

**XILINX**

# xfDNN Graph Compiler

Pass in a Network

Microcode for xDNN is Produced

xfDNN
Graph Compiler

© Copyright 2018 Xilinx

**XILINX**

# xfDNN Network Deployment



**Fused Layer Optimizations**
- Compiler can merge nodes
  - (Conv or EltWise)+Relu
  - Conv + Batch Norm
- Compiler can split nodes
  - Conv 1x1 stride 2 -> Maxpool+Conv 1x1 Stride 1

**On-Chip buffering reduces latency and increases throughput**
- xfDNN analyzes network memory needs and optimizes scheduler
  - For Fused and "One Shot" Deployment

**"One Shot" deploys entire network to FPGA**
- Optimized for fast, low latency inference
- Entire network, schedule and weights loaded only once to FPGA

XILINX.

# xfDNN Quantizer: Fast and Easy



1) Provide FP32 network and model
   - E.g., prototxt and caffemodel

2) Provide a small sample set, no labels required
   - 16 to 512 images

3) Specify desired precision
   - Quantizes to <8 bits to match Xilinx's DSP

**XILINX**

# Flexible: Multi-Network Configuration



| Application: Image Classification | Application: Speech | Application: Localization | Application: Object Detection |
|---|---|---|---|
| Model: Resnet50 | Model: Googlenet v1 | Model: FaceNet | Model: Yolo v2 |
| Framework: Caffe | Framework: MxNet | Framework: TesnorFlow | Framework: Caffe |

xfDNN

SDx Runtime

Host

PCIe

**1 FPGA Provides 4 Virtual Accelerators
For Real Time Deep Learning**

**XILINX.**

# Flexible: Bring Your own IP!

Custom Application

| Application: Localization | Application: Object Detection |
|---|---|
| Model: FaceNet | Model: Yolo v2 |
| Framework: TensorFlow | Framework: Caffe |
| xfDNN | |

SDx Runtime

Host

PCIe

xDNN PE

Custom

Platform Infrastructure

**Integrate Custom Applications Directly
with xDNN Processing Engines**

XILINX

# X + ML

XILINX.

# X + ML Focus Applications Summary

Smart City / Cloud Surveillance
• 10x Lower Latency



High Resolution Imaging
• 1.6x Faster



Security / Anomaly / Malware
• 5x Faster

XILINX

# ML Suite + Deephi

# Integrated Xilinx-Deephi Roadmap

Edge/Embedded | Cloud/DC

**Models**

20+ pruned / customized / basic models

Coming to ML Suite at XDF

Deephi Pruning

**Software Stack**

| Deephi Quantizer | | xfDNN Quantizer |
|---|---|---|
| Deephi Compiler | SDSoC \| SDAccel | xfDNN Compiler |
| Deephi Runtime | | xfDNN Runtime |

**FPGA IP**

| Deephi DPU | Deephi LSTM | xDNN |

**Platforms**

Z7020 Board    Z7020 SOM    ZU2/3 SOM    ZU2/3 Card

NIMBIX

Alibaba Cloud aliyun.com

HUAWEI

Xilinx U200, U250, U280

ZU9 Card    ZCU102    ZCU104    Ultra96

XILINX

# Xilinx Pruning Overview

## Deep compression
### Makes algorithm smaller and lighter


before pruning — after pruning — pruning synapses — pruning neurons

**Highlight**

| 1/3 | 1/10 | 1/10 | 3 X |
|-----|------|------|-----|
| Weight number | Bandwidth load | Model size | Performance |

**Compression efficiency**

Deep Compression Tool can achieve significant compression on **CNN** and **RNN**

**Accuracy**

Algorithm can be **compressed 7 times without losing accuracy** under SSD object detection framework

XILINX

# Pruning Example - SSD



Chart 1 (bar chart) — operations (G) and mAP (%):

| Index | operations (G) | mAP (%) |
|---|---|---|
| 1 | 117 | 61.5 |
| 2 | 57 | 63.4 |
| 3 | 37 | 63.5 |
| 4 | 27 | 63.4 |
| 5 | 23 | 62.4 |
| 6 | 19 | 62 |
| 7 | 17 | 61.5 |
| 8 | 15.6 | 61.1 |
| 9 | 14.6 | 61 |
| 10 | 13.6 | 60.8 |
| 11 | 12.2 | 59.2 |
| 12 | 11.6 | 60.4 |

Chart 2 (line chart) — Speedup:

| OPS | FPS |
|---|---|
| 117G | 18 |
| 19G | 71 |
| 11.6G | 103 |

XILINX

# Supported DNN (Deep Neural Network) by Applications



| Application | NTT Request | Function | Algorithm |
|---|---|---|---|
| Face | | Face detection | SSD, Densebox |
| | | Landmark Localization | Coordinates Regression |
| | | Face recognition | ResNet + Triplet / A-softmax Loss |
| | | Face attributes recognition | Classification and regression |
| Pedestrian | 1 | Pedestrian Detection (Crowd Volume) | SSD |
| | | Pose Estimation | Coordinates Regression |
| | | Person Re-identification | ResNet + Loss Fusion |
| Video Analytics | 1 | Object detection | SSD, RefineDet |
| | | Pedestrian Attributes Recognition | GoogleNet |
| | | Car Attributes Recognition | GoogleNet |
| | 1 | Car Logo Detection | DenseBox |
| | 1 | Car Logo Recognition | GoogleNet + Loss Fusion |
| | 1 | License Plate Detection | Modified DenseBox |
| | 1 | License Plate Recognition | GoogleNet + Multi-task Learning |
| ADAS/AD | | Object Detection | SSD, YOLOv2, YOLOv3 |
| | | 3D Car Detection | F-PointNet, AVOD-FPN |
| | | Lane Detection | VPGNet |
| | | Traffic Sign Detection | Modified SSD |
| | | Semantic Segmentation | FPN |
| | | Drivable Space Detection | MobilenetV2-FPN |
| | | Multi-task (Detection+Segmentation) | Deephi |

# ML Suite Performance Roadmap

ACAP

7x Performance Improvement

Model Compression,
Low and Mixed precision

**Img/Sec** — GoogleNet v1, Batch < 2

- 9000 — V100, Batch = 128
- 8000
- 7000
- 6000
- 5000
- 4000 — VCU1526 (VU13P)
- 3000 — VCU1525 (VU9P), VCU155x (VU37P)
- 2000
- 1000 — VCU1525 (VU9P), P4, V100
- K80 B8, Xeon

Legend:
- xDNNv2, INT8
- xDNNv3, INT8
- Versal AI Core, INT8

Timeline: Today, Mar'18, Jun'18, Sep'18, Dec'18, Mar'19, Jun'19

Quarters: Q4FY18, Q1FY19, Q2FY19, Q3FY19, Q4FY19, Q1FY20, Q2FY20

CPU: https://mxnet.incubator.apache.org/faq/perf.html
Nvidia: https://images.nvidia.com/content/pdf/inference-technical-overview.pdf
P4 = int8, v100 = fp16

© Copyright 2018 Xilinx

XILINX

# Visit Xilinx.com/ML for more information

[https://www.xilinx.com/applications/megatrends/machine-learning.html](https://www.xilinx.com/applications/megatrends/machine-learning.html)

# Adaptable.
# Intelligent.

**XILINX.**