# CEVA Deep Neural Network
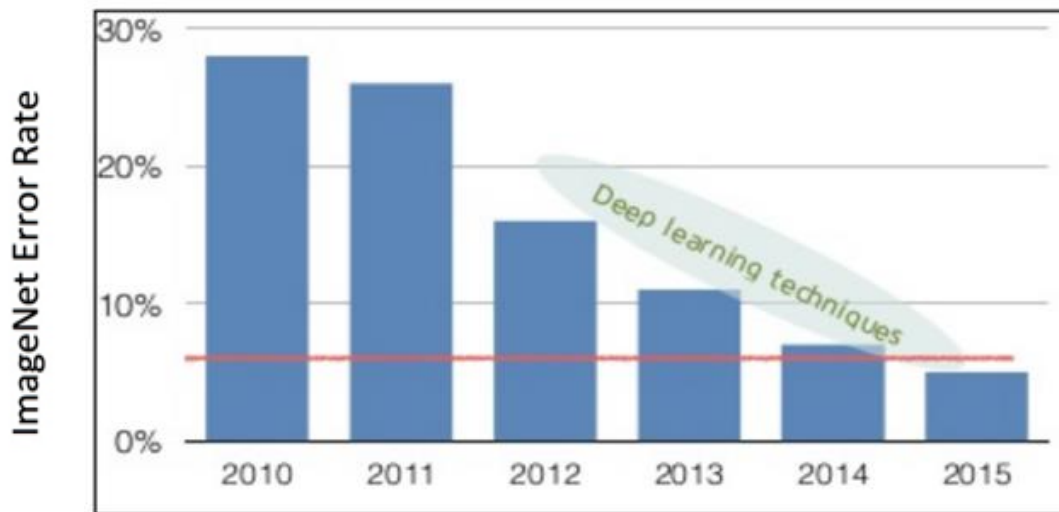
Liran Bar, Director of Product Marketing (Liran.Bar@ceva-dsp.com)

July 2016 – Please contact us for the update material

www.ceva-dsp.com

# Deep Learning Performance Improvements

# The Neural Networks Challenge
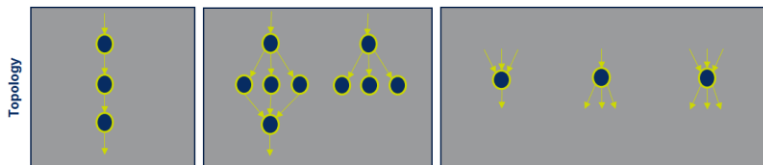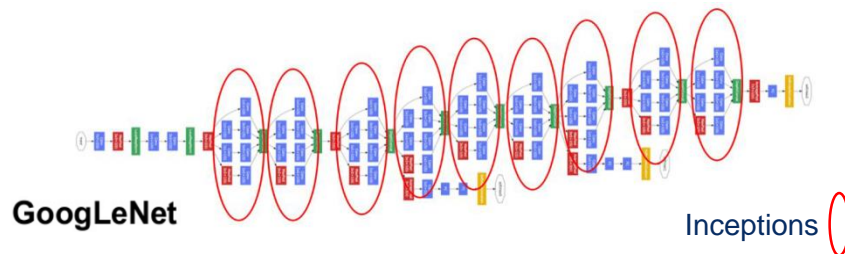
▶ Various training **frameworks**: Caffe, TensorFlow, Torch etc.

▶ Various **layers**: Convolutional, Normalization, Pooling etc.

▶ Various network **topologies**



▶ Need to deal with network inside a network



GoogLeNet

Inceptions

▶ Need to have optimized solution for variable sizes of ROI

# Leading Deep Learning Frameworks

## Caffe

- A well-known and widely used machine-vision library
- Implementation of fast convolutional nets to C and C++
- Made with expression, speed, and modularity in mind
- Used by researchers, academy, and commercial companies

## Google TensorFlow

- Relatively new alternative to Caffe, supported and promoted by Google
- Scalable to work both for research and commercial purpose without making any changes
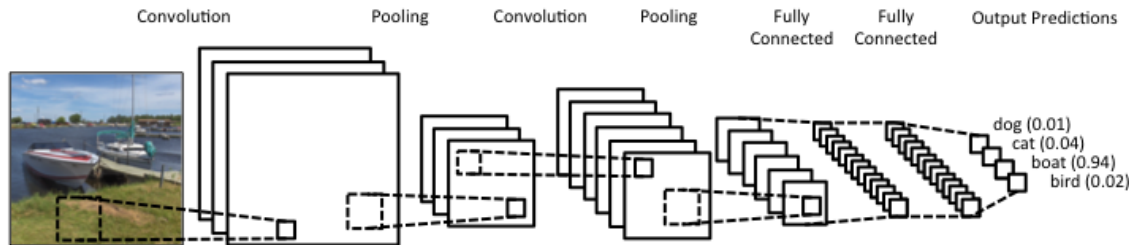- A software library for numerical computation using data flow graphs

# Leading Neural Network Layers

CDNN2

CEVA®

▶ Convolutional

▶ Normalization

▶ Pooling  (Average and Max)

▶ Fully Connected

▶ Activation (ReLU, Parametric ReLU, TanH, Sigmoid)

▶ Deconvolution

▶ Concatenation

▶ Upsample

▶ Argmax

▶ Softmax

Flexible embedded solution is required to cope with the evolving and leading neural network layers



Convolution    Pooling    Convolution    Pooling    Fully Connected    Fully Connected    Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

# Deep Learning Topologies



| Linear Networks | Multiple Layers Per Level | Multiple-Input-Multiple-Output |
|---|---|---|
| | | (a)   (b)   (c) |
| AlexNet<br>VGG-19<br>VGG-16<br>VGG_S | GoogLeNet | GoogLeNet<br>SegNet<br>ResNet |

**Topology**

**Networks**

# Neural Network Embedded Challenges

Implementing a deep neural network in embedded systems is a challenging task!

▶ Very high bandwidth consuming and computing bottleneck
- ▶ Data transfer in/out the DDR – Input/output maps
- ▶ Convolution and Fully Connected data weights from DDR
- ▶ Processing multi ROI with the same network
- ▶ Internal memory size limitation on embedded platforms

▶ Porting and optimizing neural networks to embedded system consumes time!
- ▶ Special programing knowledge (intrinsic, assembly)
- ▶ Specific experience in the embedded platform (instructions, hardware capabilities)
- ▶ Fixed HW solutions are not flexible to cope with the evolving neural networks
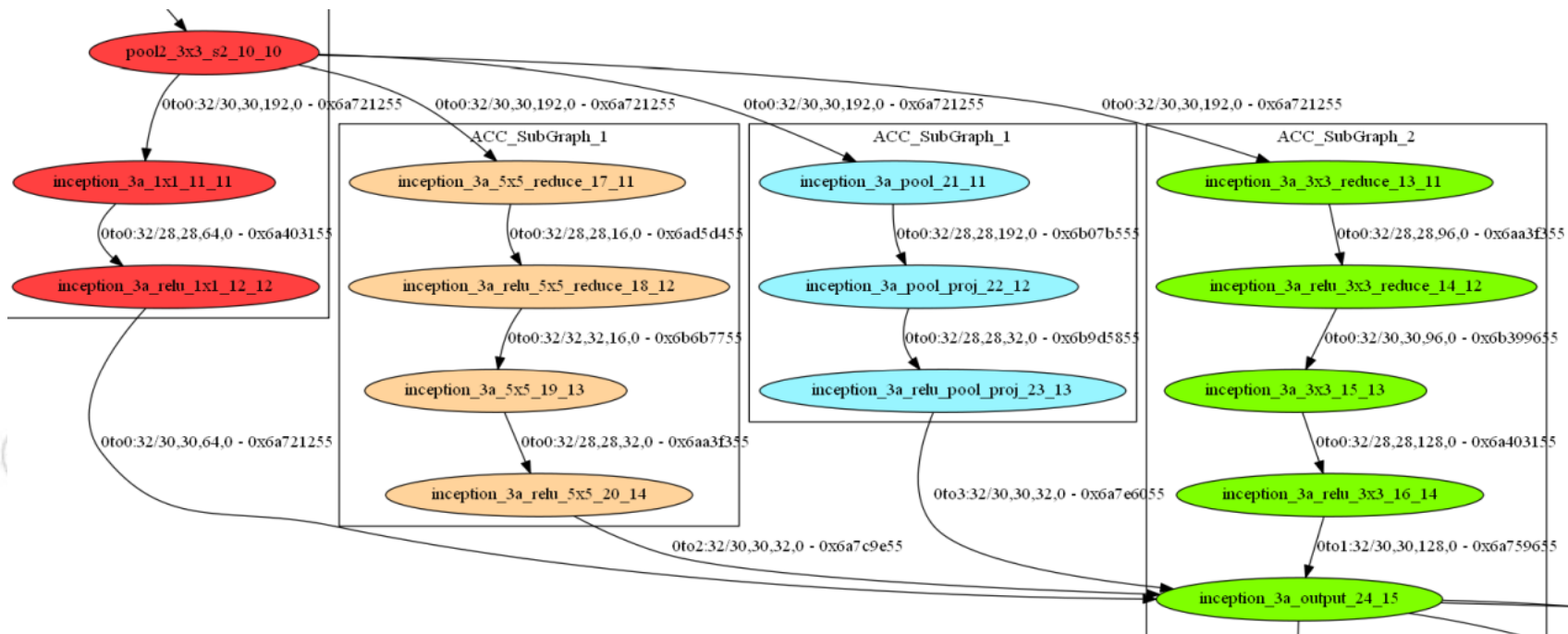
## Long "Time-To-Market"

# GoogleNet Challenge

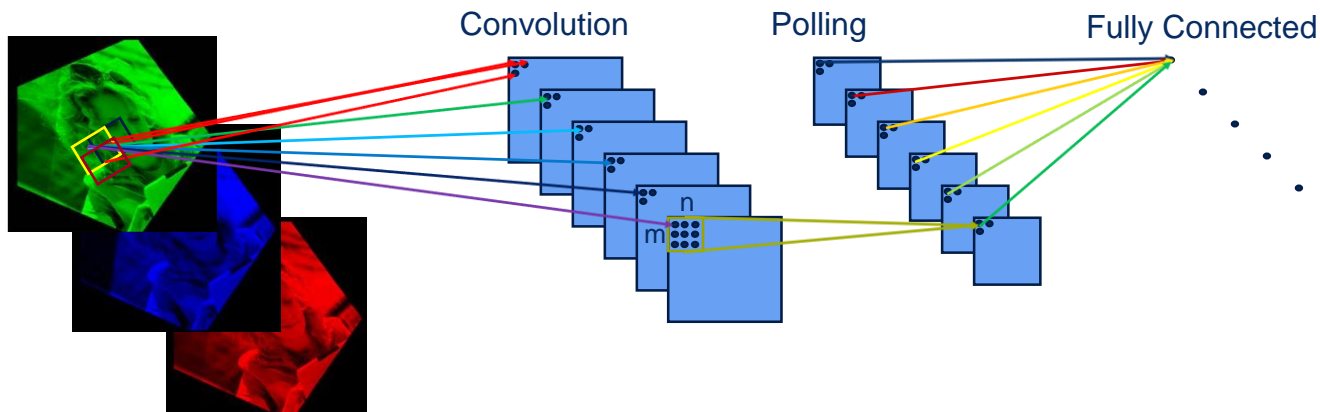▶ **How do you run a network like GoogLeNet with minimum transaction through DRR?**

# What can be done?

▶ **Reducing bandwidth**

  ▶ In the convolution layer, each output is calculated by the same inputs

    ▶ Weights matrix are shared between output results in the same map (in order not to load the weights more than once)

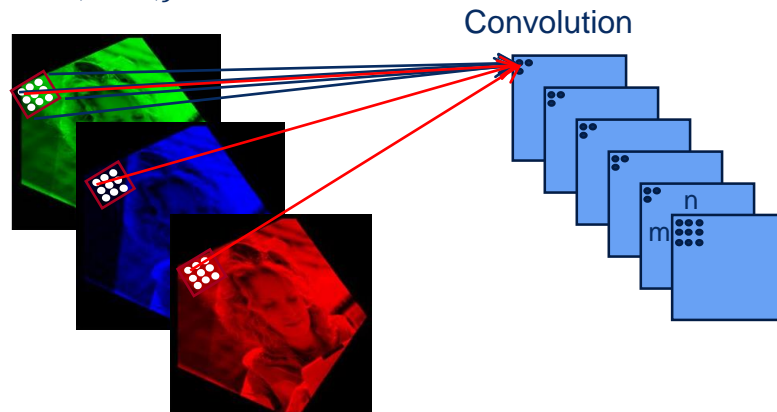    ▶ The input data can be reused to avoid useless transactions from DDR



Convolution          Polling          Fully Connected

# What can be done? – Cont.

▶ **Maximum multiply accumulate utilization**

  ▶ Differentiate between large inputs to small input maps and the number of maps from each type

  ▶ Large size maps - $X_{c,i,j}^l = \sum_{c=0}^{C} \sum_{i=0}^{H} \sum_{j=0}^{W} W_{c,m,n}^l X_{c,i+m,j+n}^{l-1}$

  ▶ Large amount of maps with small size (last layers) - $X_{c,i,j}^l = \sum_{i=0}^{H} \sum_{j=0}^{W} \sum_{c=0}^{C} W_{c,m,n}^l X_{c,i+m,j+n}^{l-1}$

Convolution

# What can be done? – Cont.

▶ **Overcome small internal memory size**

    ▶ Try to preserve the principle of "All inputs must be in the internal memory" by tile division

    ▶ Divide all input maps to identical tile sizes

Convolution

# What can be done? – Cont.

- Use compression algorithms and prior knowledge to reduce bandwidth to and from the external memory
  - Using algorithms like huffman coding
  - Work in pipeline to save BW
  - Identify when some of the calculation can be saved
  - Share data between calculations
  - Recognize when the focus should be on the weights and when it should be on the map size – network dependent
  - Compress and decompress better over time (learn from frame by frame execution)

# CEVA Deep Neural Network (CDNN2)

▶ **2<sup>nd</sup> gen SW framework support**

> ▶ Caffe and TensorFlow Frameworks
>
> ▶ Various networks*
>
> ▶ All network topologies
>
> ▶ All the leading layers
>
> ▶ Variable ROI
>
> ▶ "Push-button" conversion from pre-trained networks to optimized real-time
>
> ▶ Accelerates machine learning deployment for embedded systems
>
> ▶ Optimized for CEVA-XM4 vision DSP

(*) Including AlexNet, GoogLeNet, ResNet, SegNet, VGG, NIN and others

# CDNN2 Usage Flow

# Real-Time CDNN2 Application Flow

# AlexNet - Network Performance

▶ Network specification

  ▶ Full forward classification case image measurement (single iteration)

  ▶ 24 layers, 224x224 network input size

  ▶ 11x11, 5x5 and 3x3 convolution filters

▶ Memory bandwidth

  ▶ Pre-trained network: 253Mbytes floating point

  ▶ Post CDNN2 (optimized for CEVA-XM4): 16Mbytes fixed-point
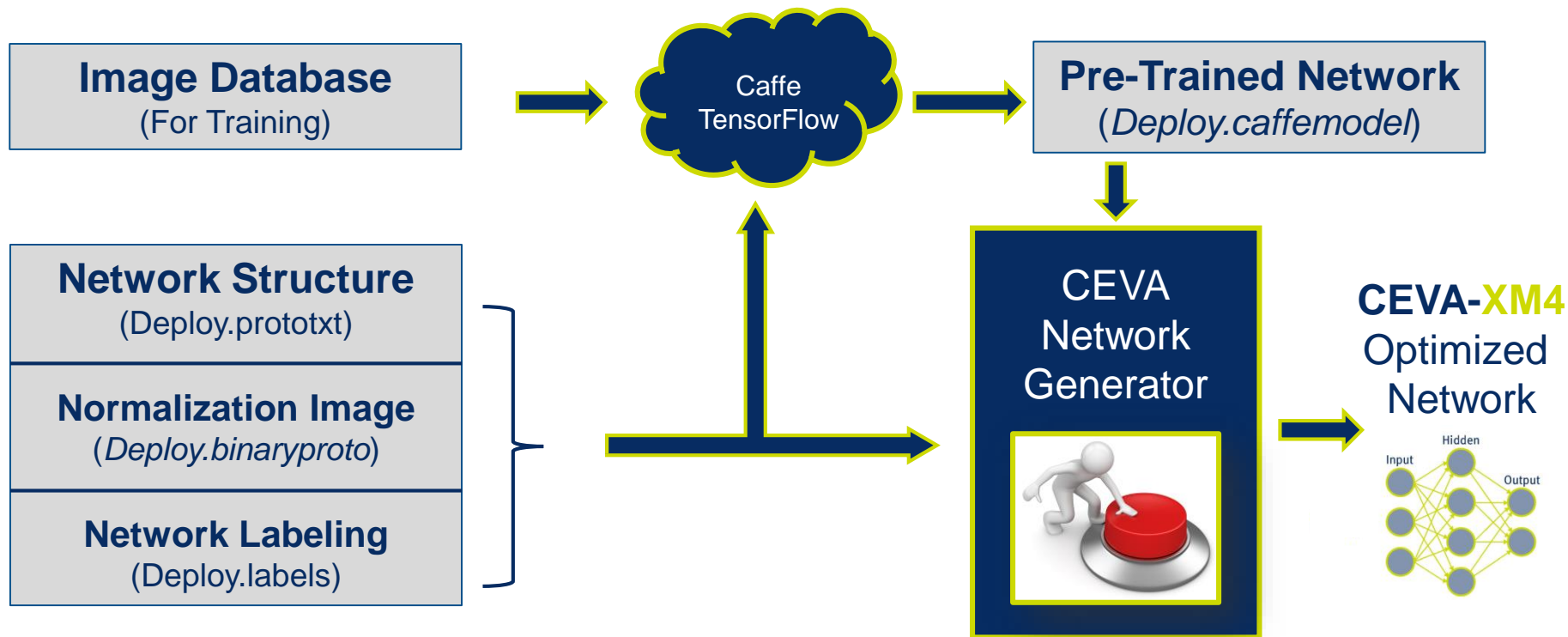
  ▶ Including weights and data

# CEVA Network Generator

# Real-Time Network Generator Demo

Live CDNN2 demo:
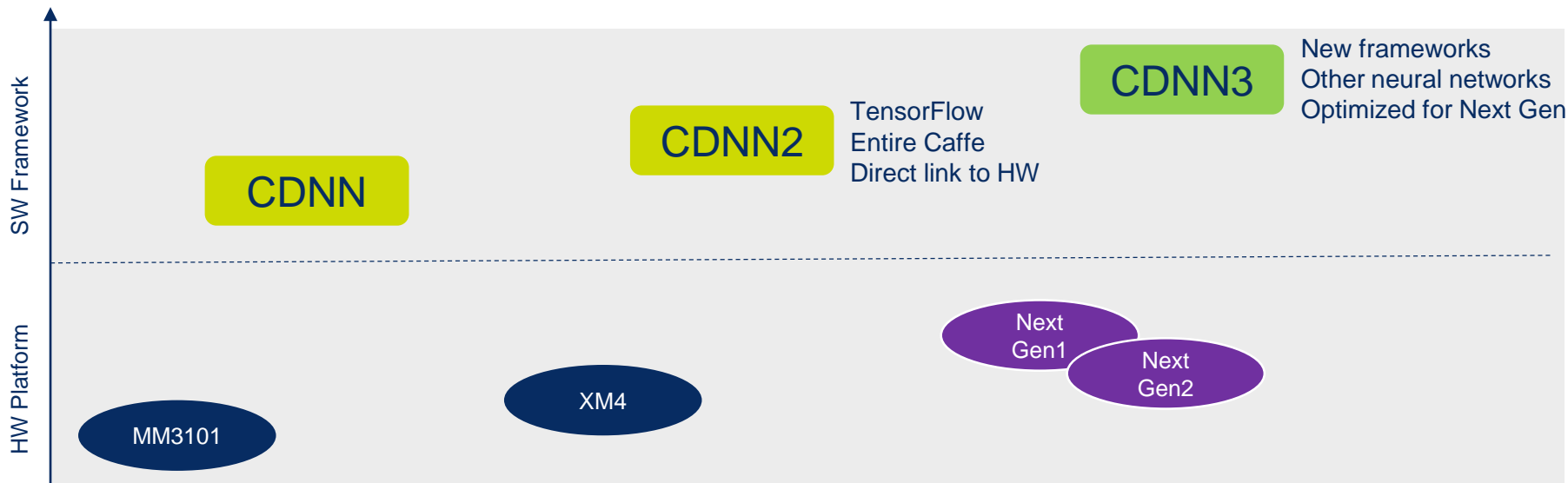https://www.youtube.com/watch?v=SXINFryLM3Q&feature=youtu.be



Downloading Age classification
Neural Network from the internet

Passing it via CEVA Network Generator and
running it on the XM4 FPGA **under 10 min !**

# Final Comments: HW + SW

▶ CEVA is at the forefront of development of Neural Network embedded platforms

▶ Normally, the HW platform is meaningless if not supported by the corresponding SW framework…



SW Framework

CDNN

CDNN2
TensorFlow
Entire Caffe
Direct link to HW

CDNN3
New frameworks
Other neural networks
Optimized for Next Gen

HW Platform

MM3101

XM4

Next Gen1

Next Gen2

CEVA®
The DSP Powerhouse

Thank You

www.ceva-dsp.com