

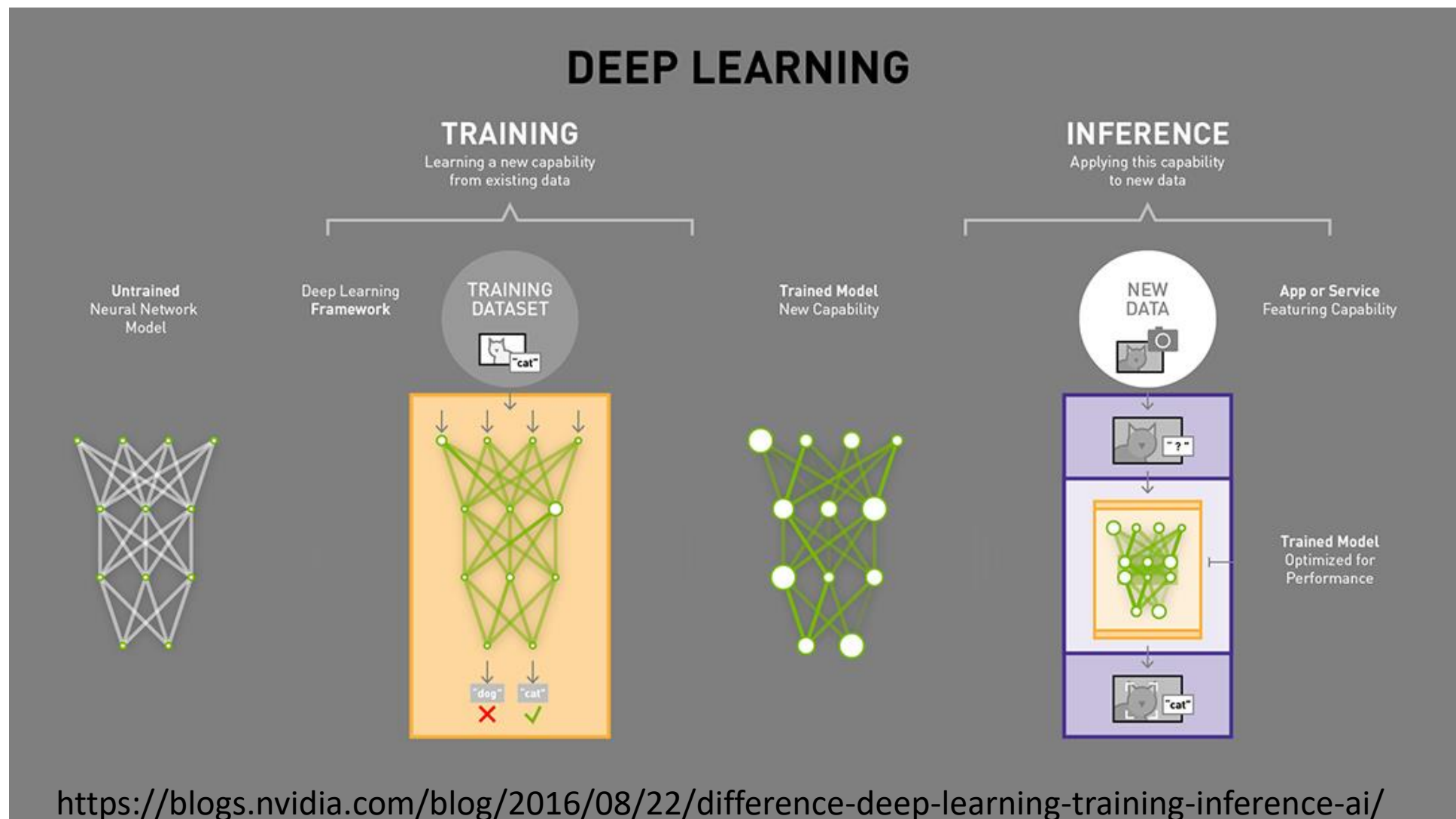
Cognitive computing at the edge

Paolo Meloni, PhD

DIEE- Università degli Studi di Cagliari

Deep Learning at the edge

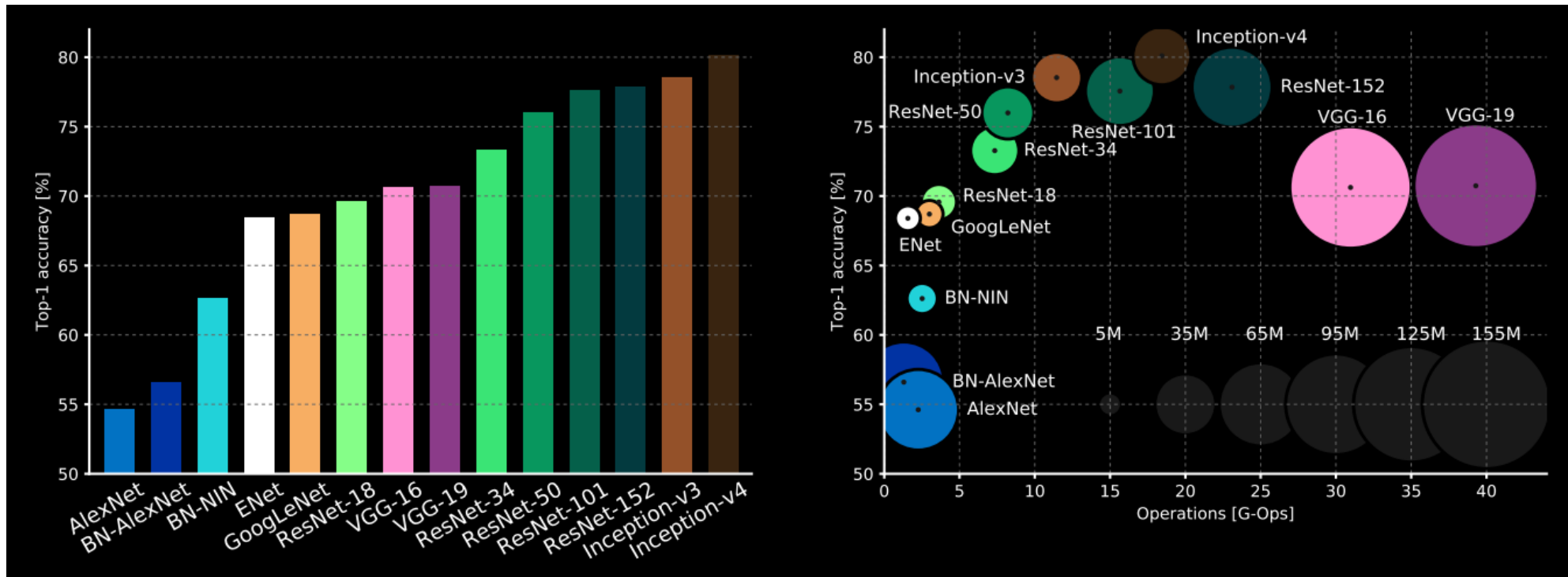
- Challenges
- Efficient implementation techniques
 - Algorithms
 - Platforms
 - Tools



Why edge computing?

- Latency
- Resilience
- Privacy

Challenges: Model size and complexity

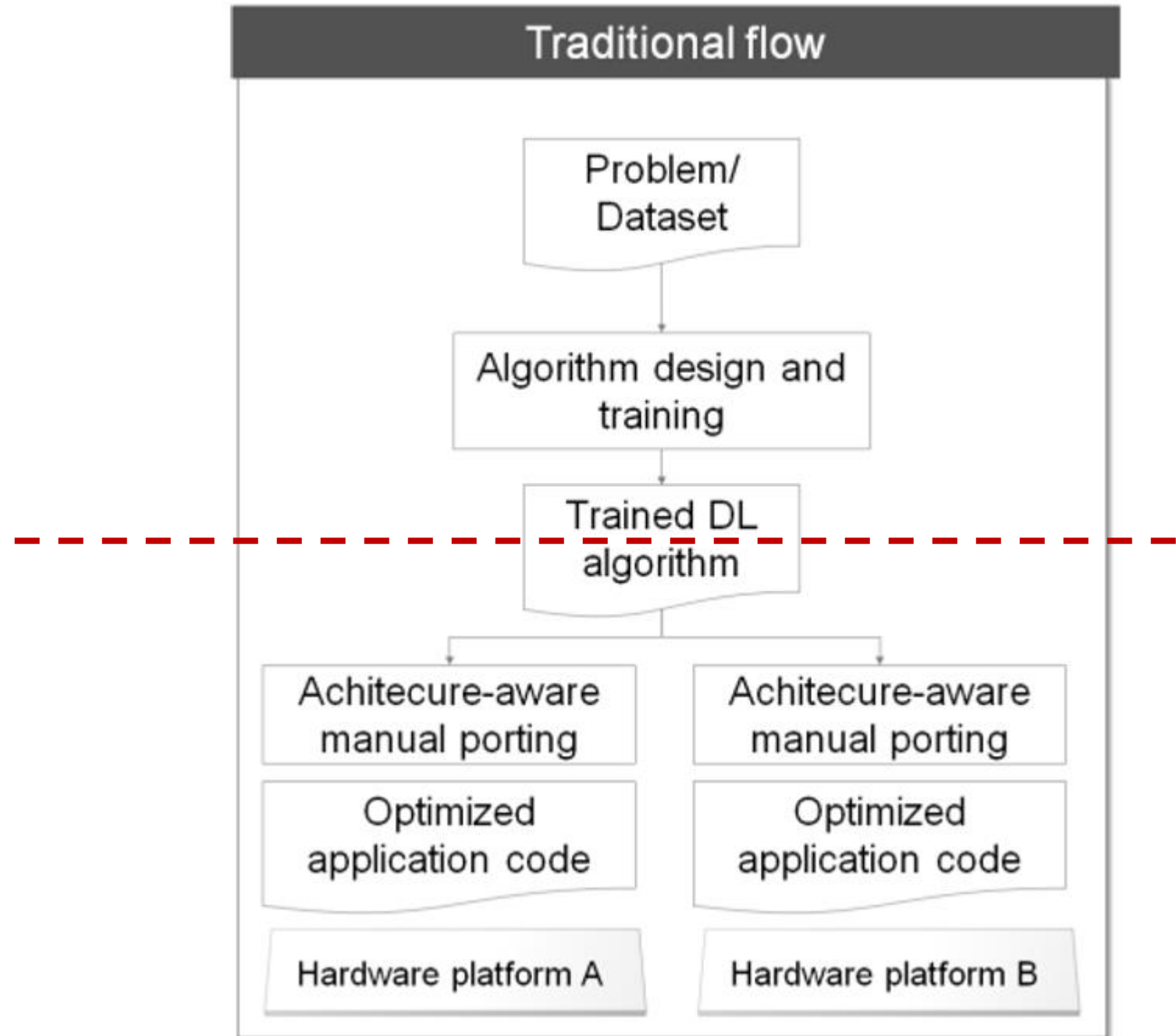


An Analysis of Deep Neural Network Models for Practical Applications A. Canziani, A. Paszke, E. Culurciello, 2016

Enormous computational and memory requirements

e.g., VGG-19: 140 million floating-point parameters to classify a single image

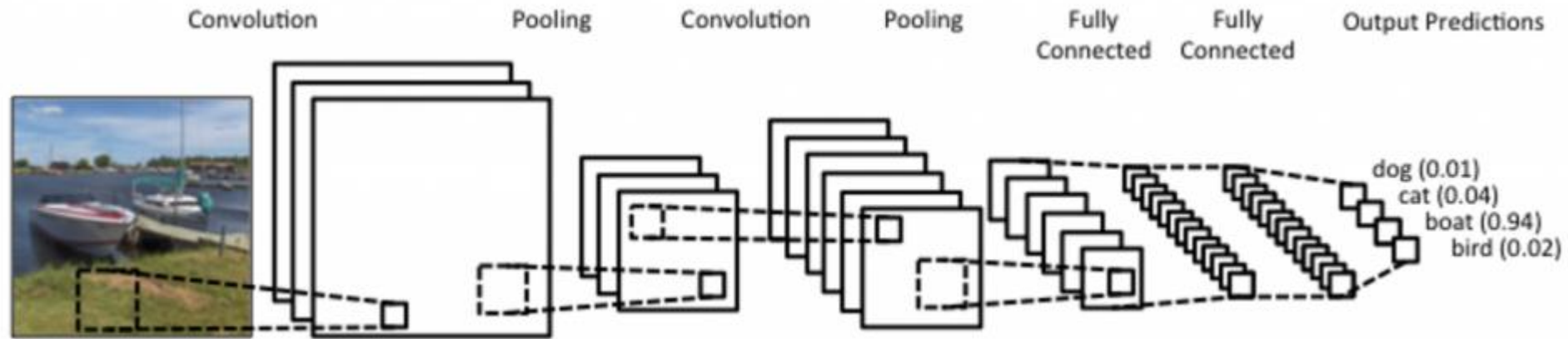
The DL dichotomy



Efficient implementation techniques

- Algorithms
- Platforms
- Tools

Algorithms: example – Convolutional Neural Network



1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

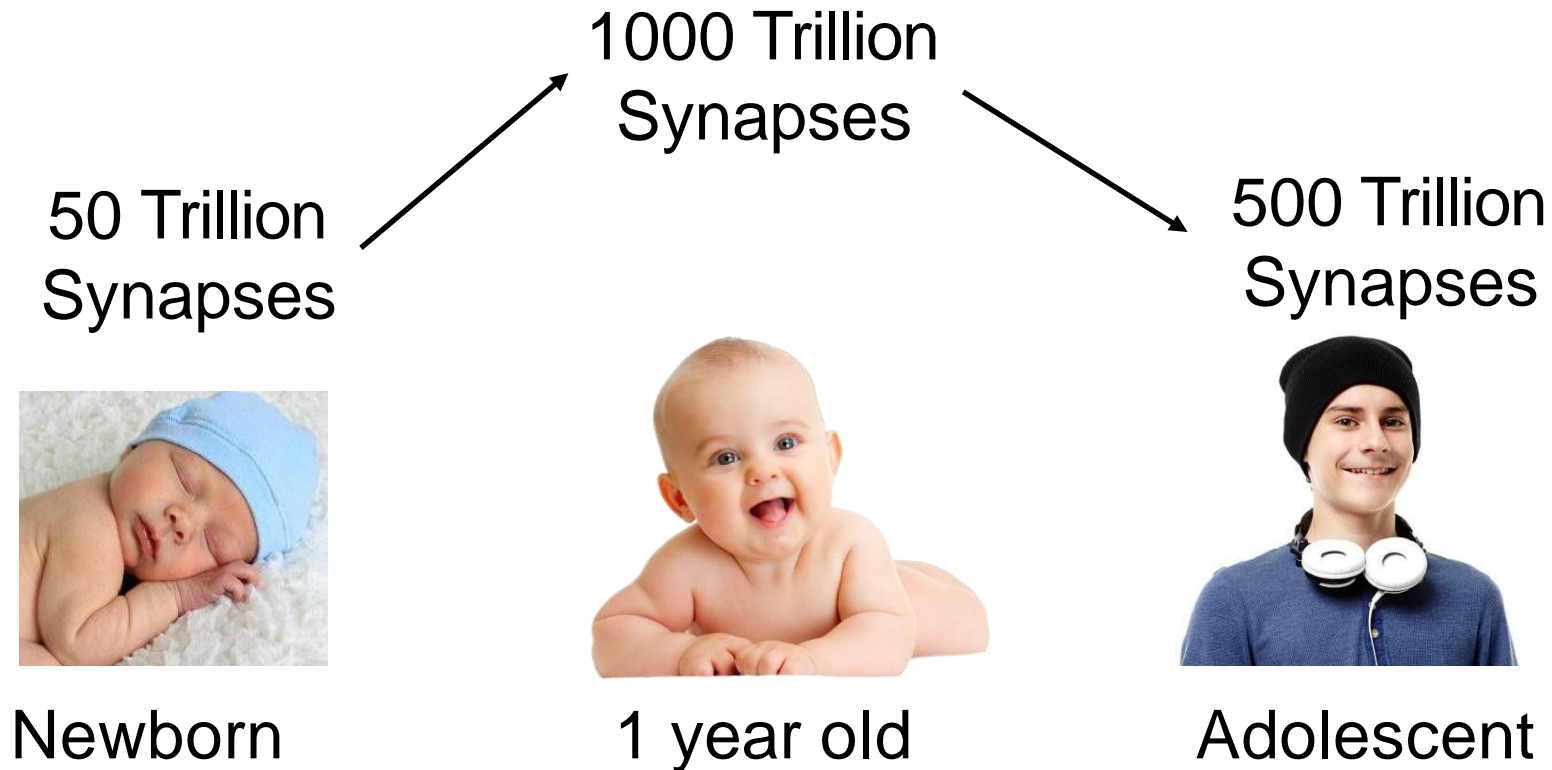
4		

Convolved Feature

Convolution with 3×3 Filter. Source:

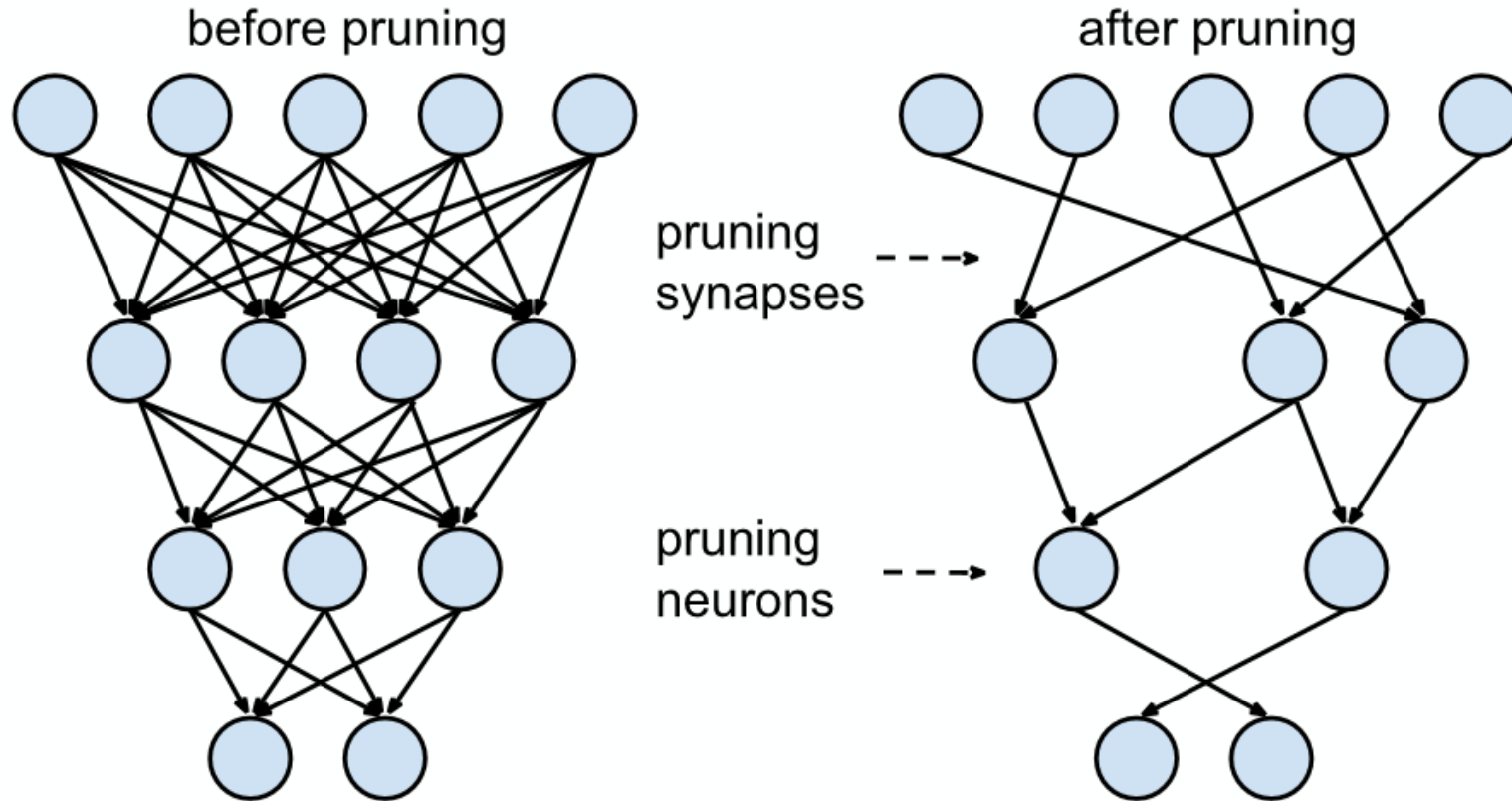
http://deeplearning.stanford.edu/wiki/index.php/Feature_extraction_using_convolution

#Synapses in Human Brain

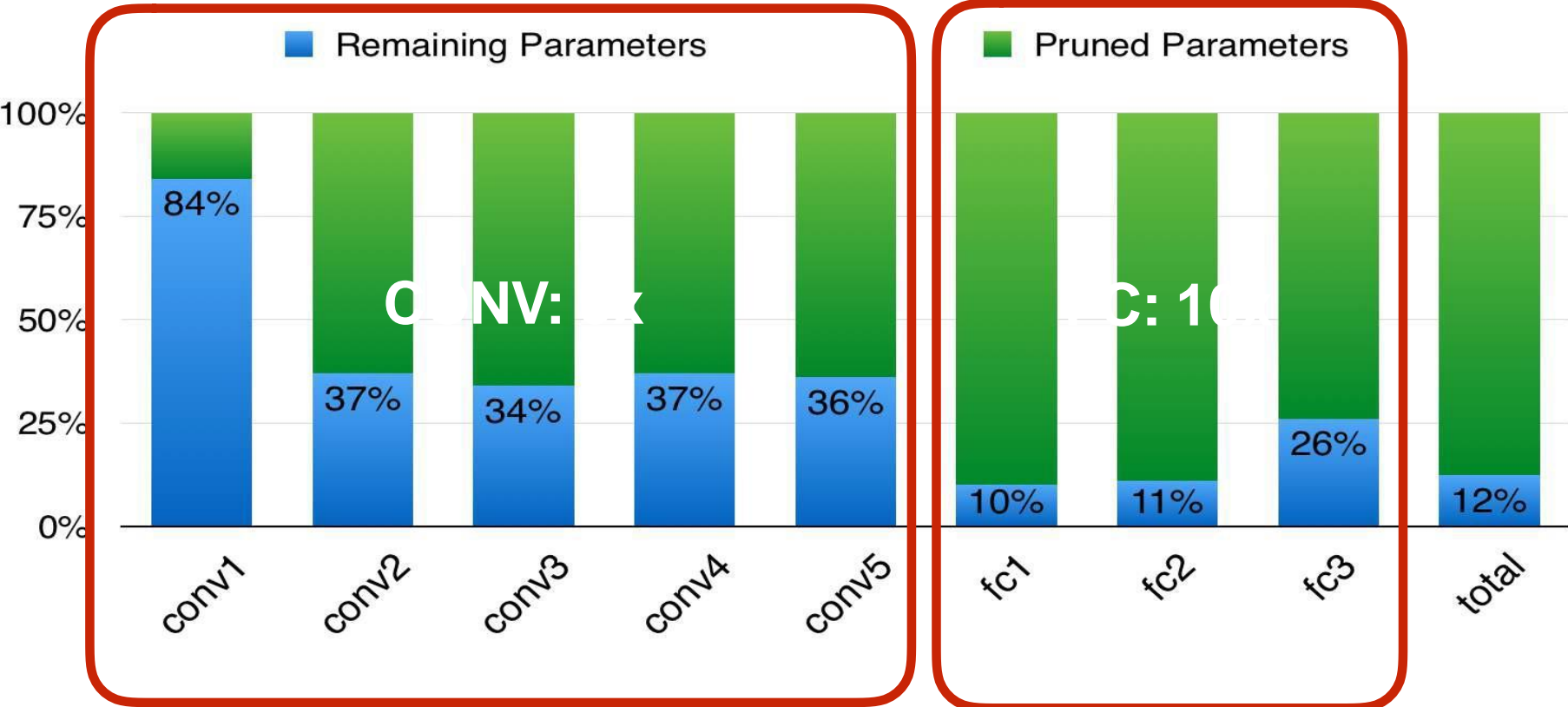


Christopher A Walsh. Peter Huttenlocher (1931-2013). *Nature*, 502(7470):172–172, 2013.

Pruning Neural Networks



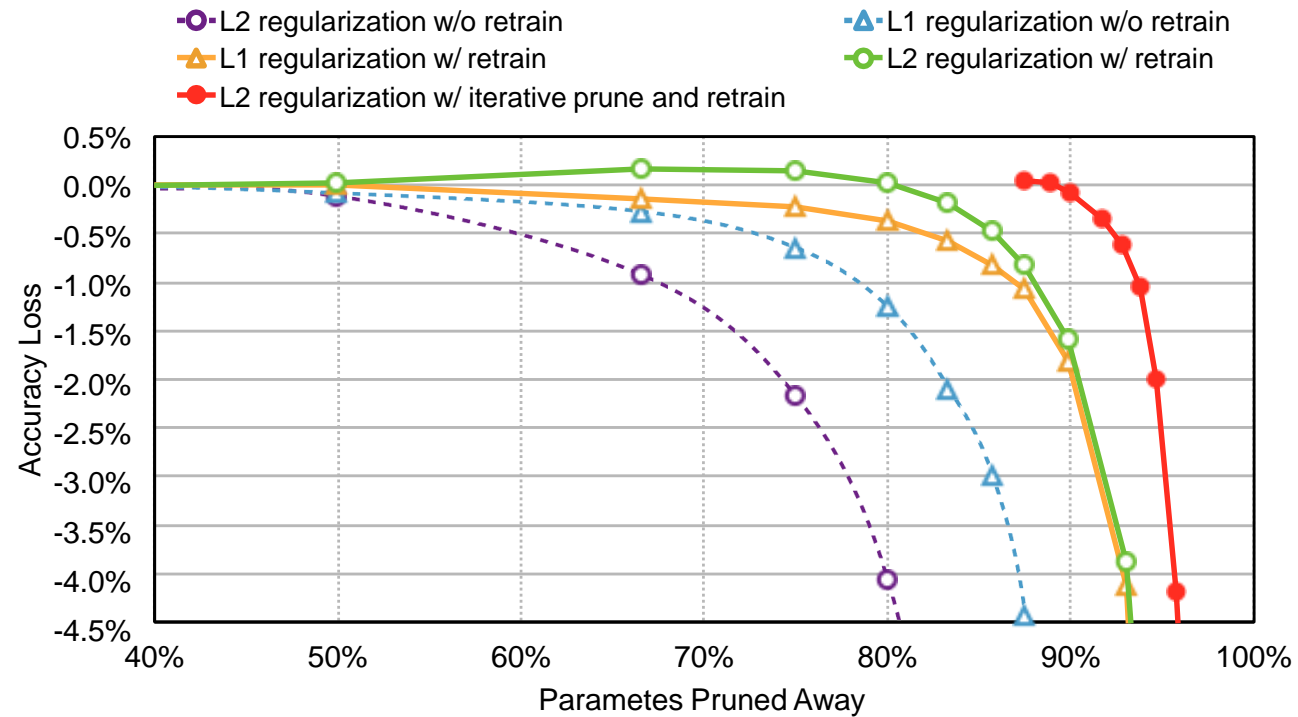
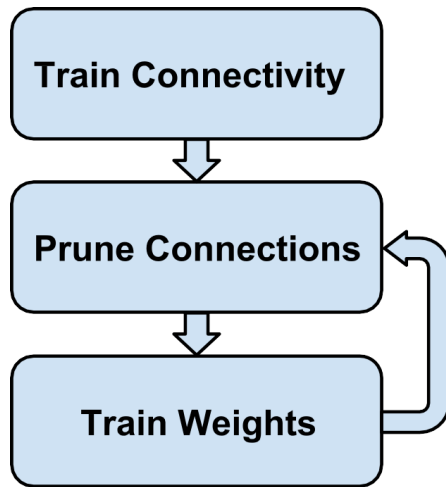
AlexNet



Han et al. Learning both Weights and Connections for Efficient Neural Networks, NIPS 2015

[Han16]

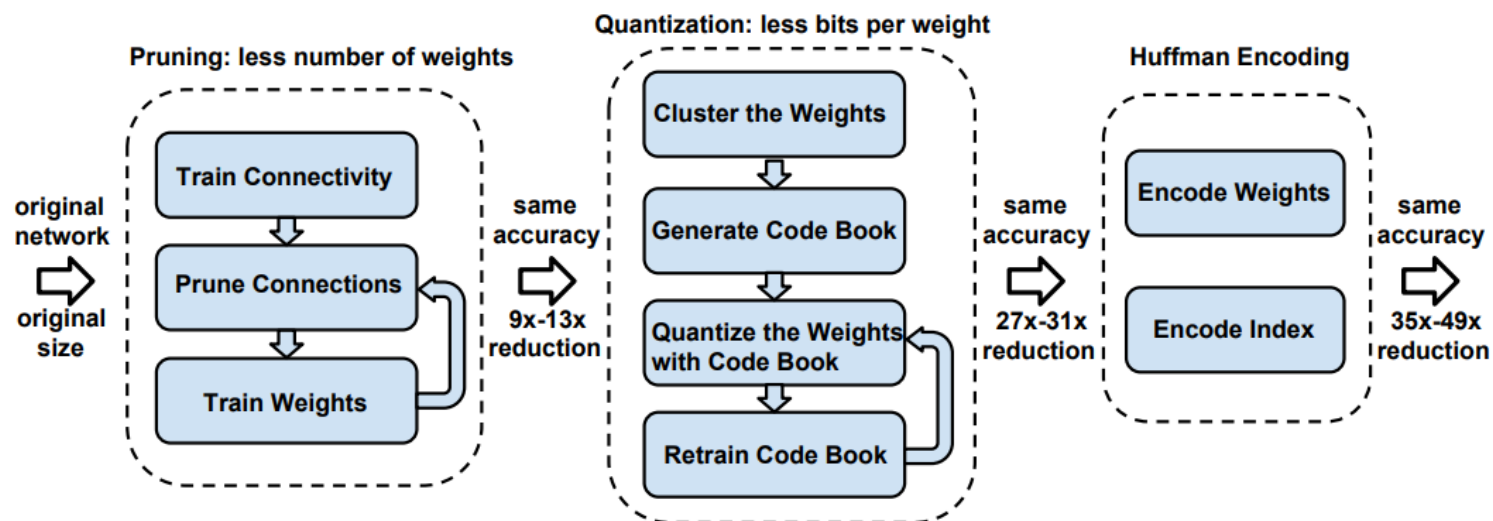
Retrain to Recover Accuracy



Han et al. Learning both Weights and Connections for Efficient Neural Networks, NIPS 2015

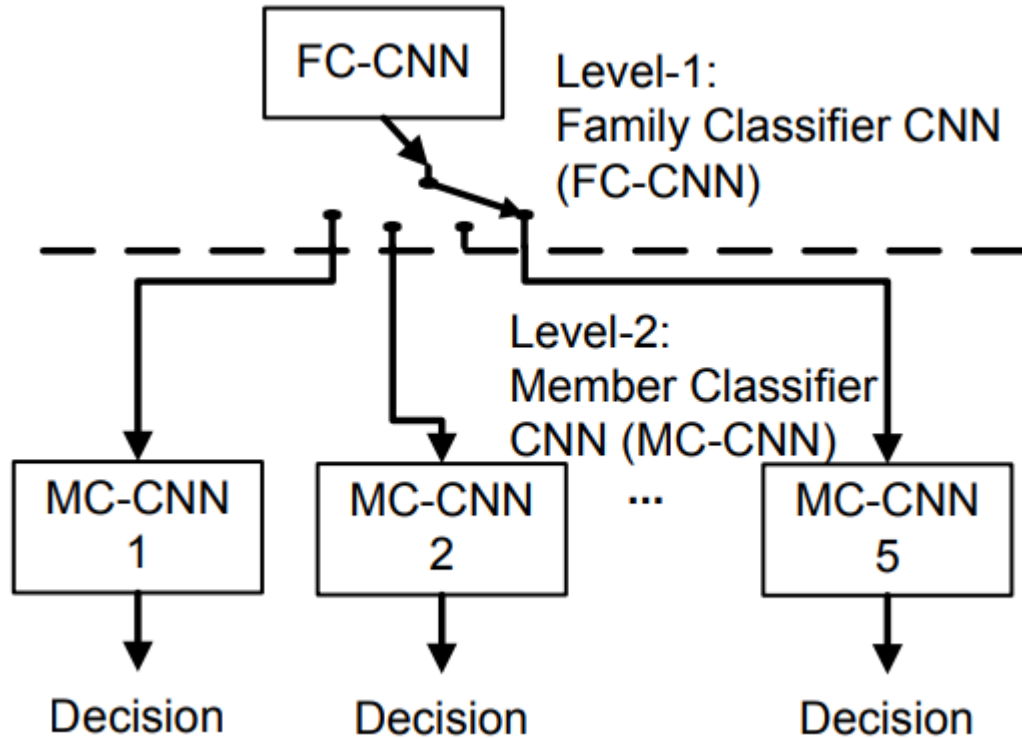
[Han16]

Quantization and Compression

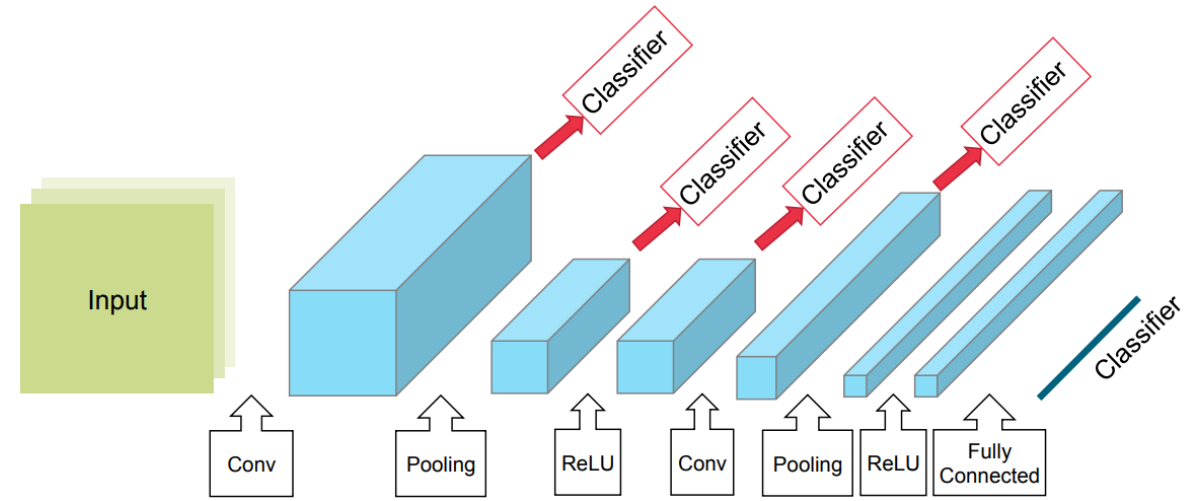


Network	Top-1 Error	Top-5 Error	Parameters	Compress Rate
LeNet-300-100 Ref	1.64%	-	1070 KB	
LeNet-300-100 Compressed	1.58%	-	27 KB	40×
LeNet-5 Ref	0.80%	-	1720 KB	
LeNet-5 Compressed	0.74%	-	44 KB	39×
AlexNet Ref	42.78%	19.73%	240 MB	
AlexNet Compressed	42.78%	19.70%	6.9 MB	35×
VGG-16 Ref	31.50%	11.32%	552 MB	
VGG-16 Compressed	31.17%	10.91%	11.3 MB	49×

Hierarchical CNN



- 5+1 different CNNs for classification
- lower complexity than one-vs-all classifier

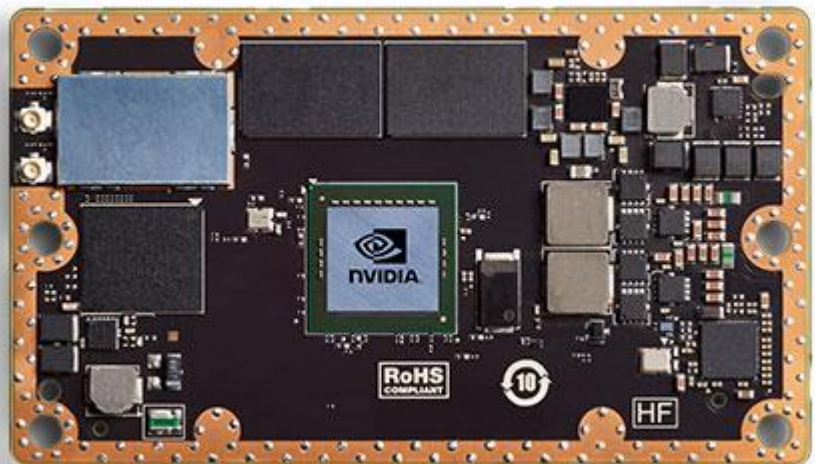


Family 1	
Family 2	
Family 3	
Family 4	
Family 5	

Platforms

- GPU
- FPGA
- Specialized processing elements

GPUs – NVIDIA



- >90% utilization of processing elements
- Good floating point support
- Good tool support CUDNN
- Good performance on vanilla CNNs

Network: AlexNet	Batch Size	Tegra X1 (FP32)	Tegra X1 (FP16)	Core i7 6700K (FP32)
Inference Performance	1	47 img/sec	67 img/sec	62 img/sec
Power		5.5 W	5.1 W	49.7 W
Performance/Watt		8.6 img/sec/W	13.1 img/sec/W	1.3 img/sec/W

https://www.nvidia.com/content/tegra/embedded-systems/pdf/jetson_tx1_whitepaper.pdf

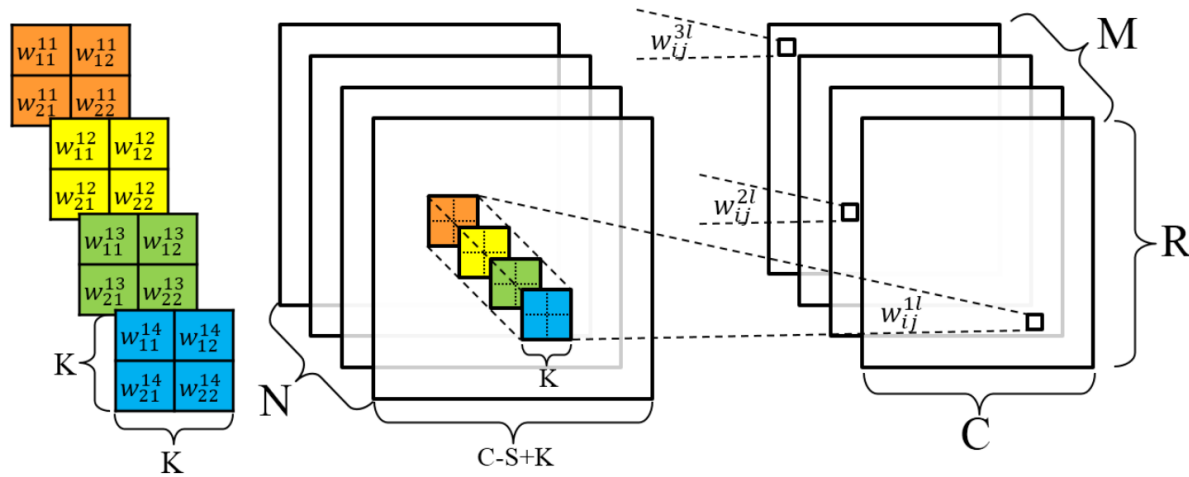
GPU - MALI

- CaffeNet (similar to AlexNet @around 10 FPS on mobile phone)

Table 1: inference speed of CaffeNet on ILSVRC 2012 and FOOD 101 datasets (Batch 12)

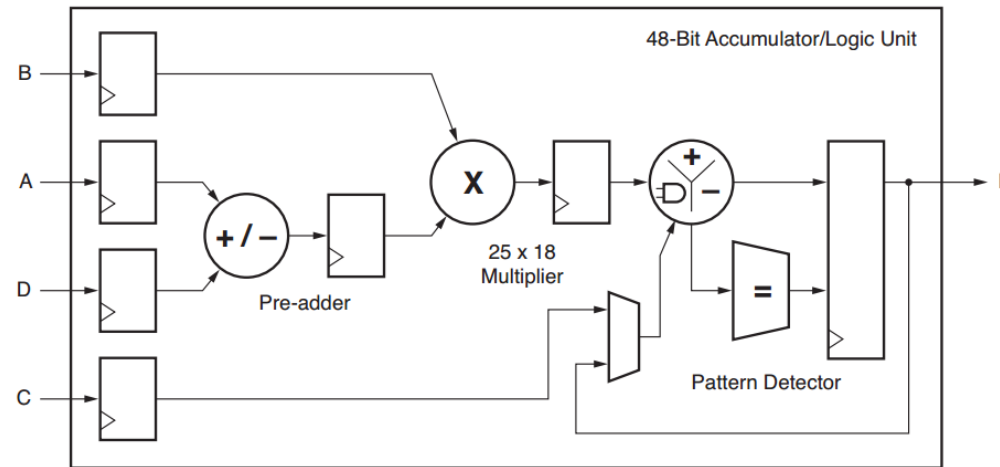
Image Dataset	Mean, inference time
Xiaomi Redmi Note 4: Mediatek MT6797 Helio X20 (ARM/Mali-T880 MP4)	
ILSVRC 2012	250.7 ms
FOOD 101	250.2 ms
Samsung S7 Edge: Exynos 8890 Octa (ARM/MALI-T880 MP 12)	
ILSVRC 2012	110 ms
FOOD 101	110 ms

Convolutional Neural Networks on FPGAs – CONV layer to DSPs



```

for i in range(0, nb_out_feat):
    for j in range(0, nb_in_feat):
        for s in range(0, feat_height):
            for t in range(0, feat_width):
                Conv= Conv (X[j,s,t] * Wt[i,j])
                //sum of product
                Y[i,s,t] += conv
    
```

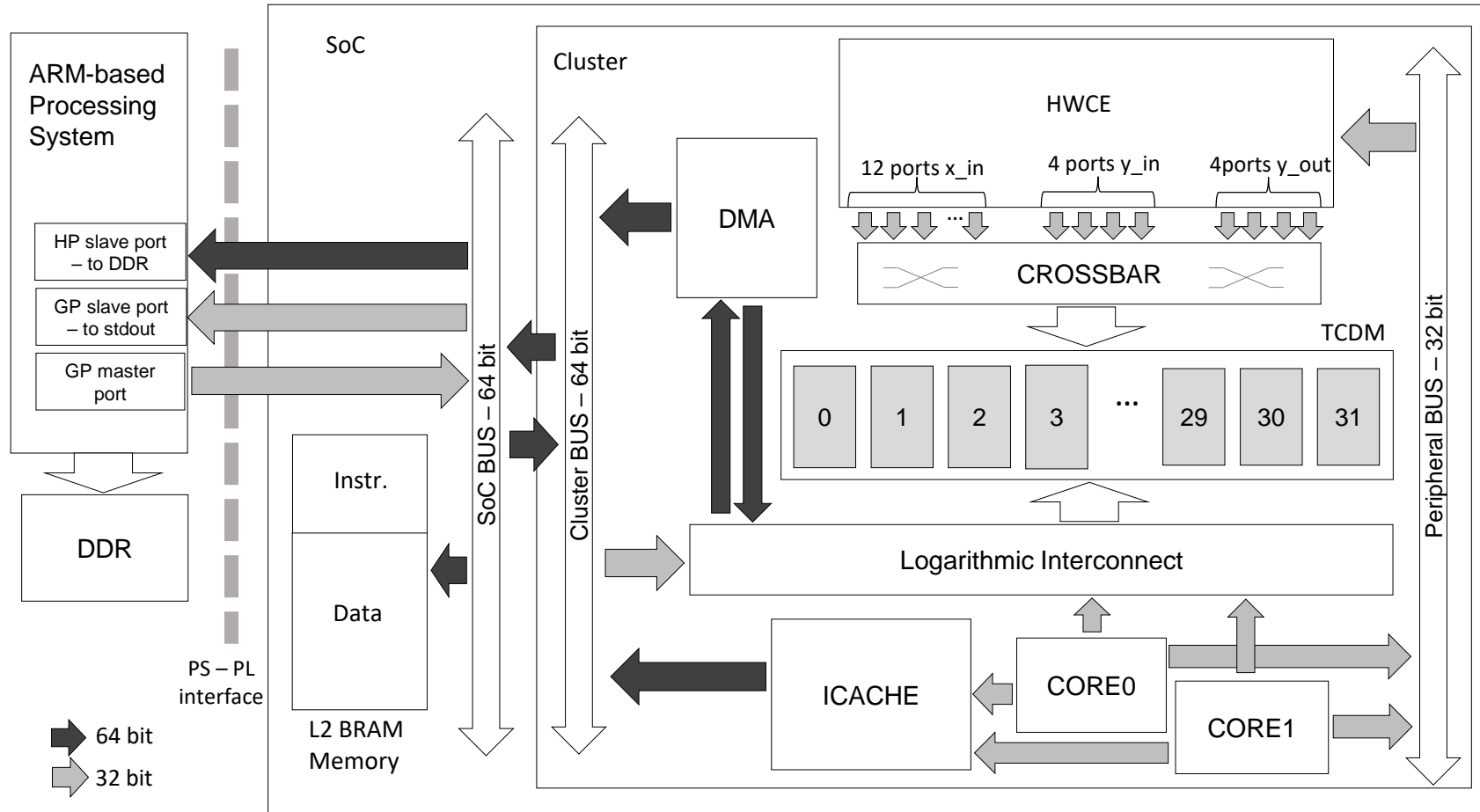


UG479_e1_21_032111

- Multiply-and-accumulate intensive
- High level of parallelism
- Maps well on FPGA MAC primitives (Xilinx DSP48 slices)

[Zhang et al., 2015]

CNN acceleration on FPGAs: Neuraghe

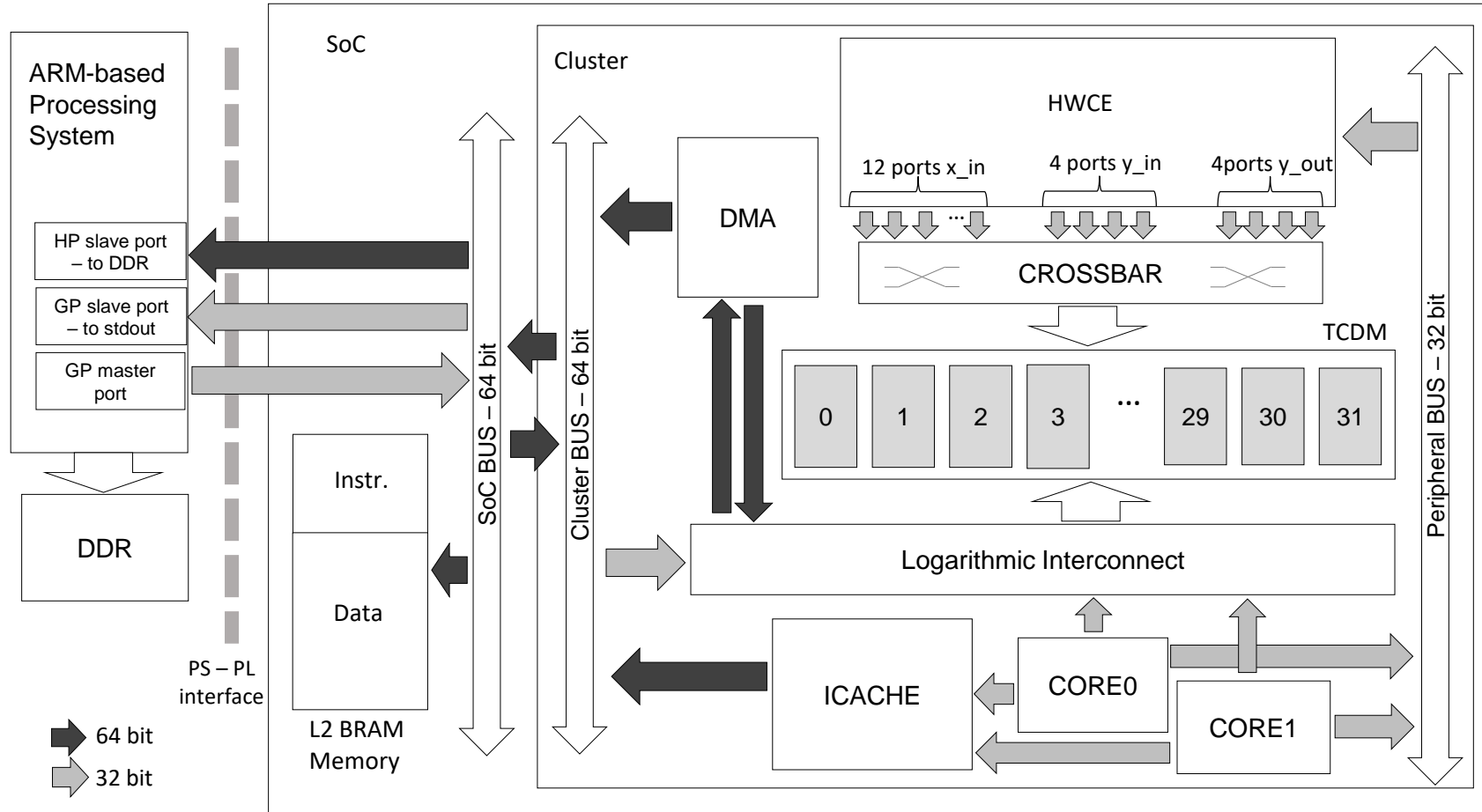


P. Meloni, G. Deriu, F. Conti, I. Loi, L. Raffo and L. Benini, "A high-efficiency runtime reconfigurable IP for CNN acceleration on a mid-range all-programmable SoC," *2016 International Conference on ReConfigurable Computing and FPGAs (ReConFig)*, Cancun, 2016, pp. 1-8.

Design principles

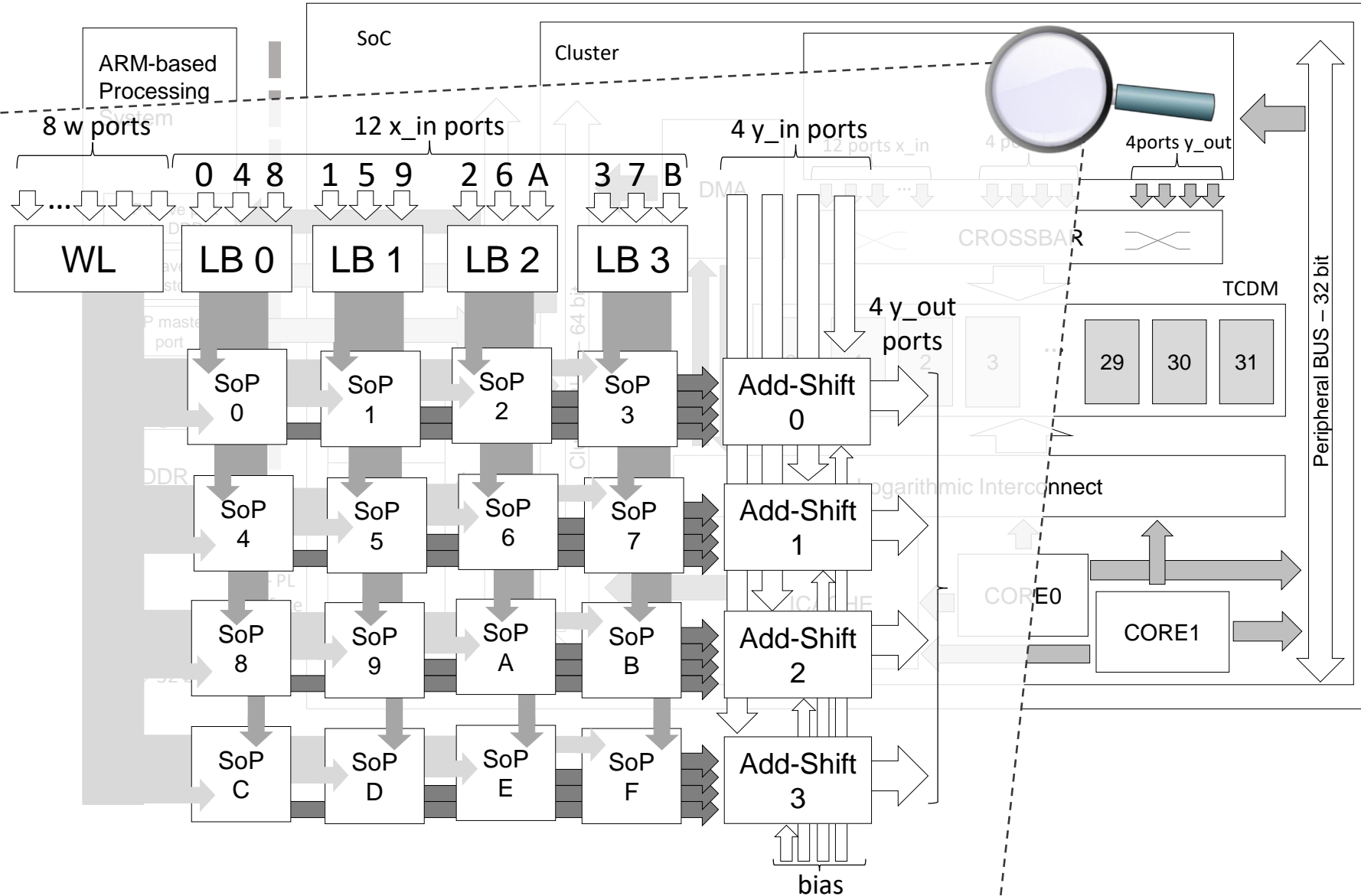
- Exploit parallelism
 - Use as many DSP as possible (Reference SoC – Zynq Z7045)
- Support dynamic reconfiguration
 - Adapt to multiple conv layers
 - Different kernel sizes
 - Different strides
- Reduce I/O bottleneck
 - Careful convolution scheduling
- All-programmable SoCs
 - Exploit ARM for housekeeping and other CNN layers

Architectural template

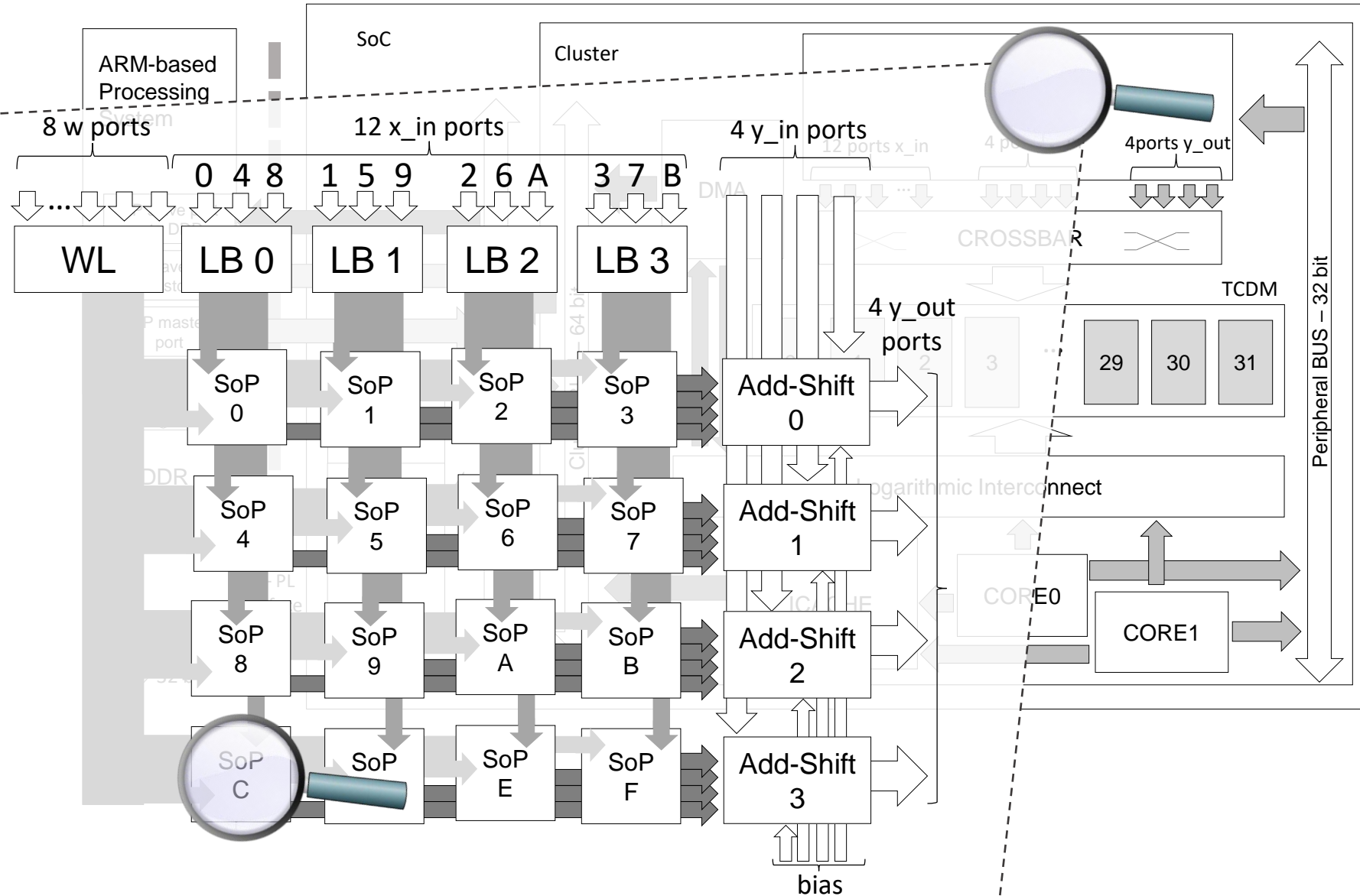


F. Conti et al. PULP: A Ultra-Low Power Parallel Accelerator for Energy-Efficient and Flexible Embedded Vision.
Journal of Signal Processing Systems, 2015.

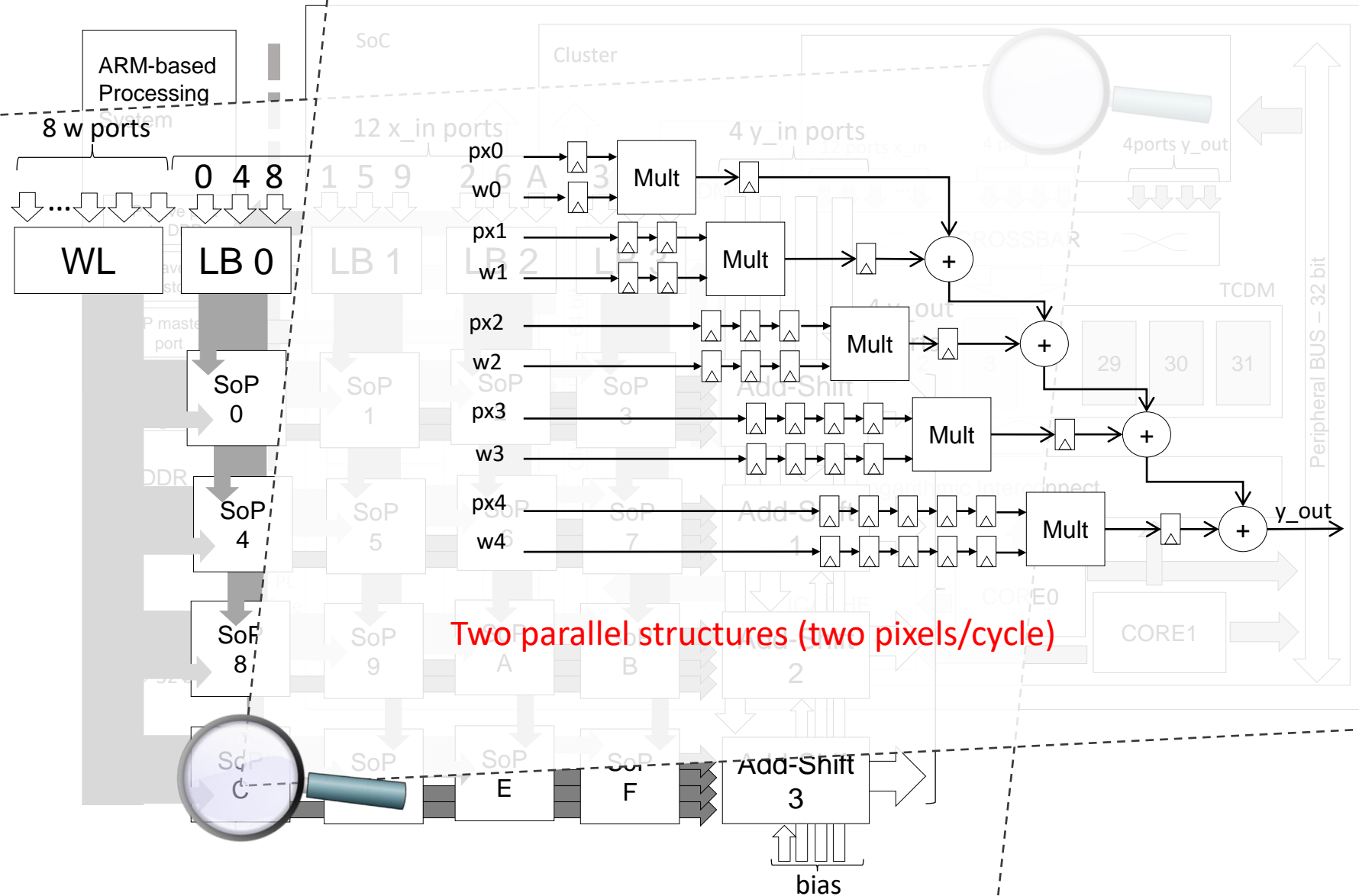
Architectural template



Architectural template



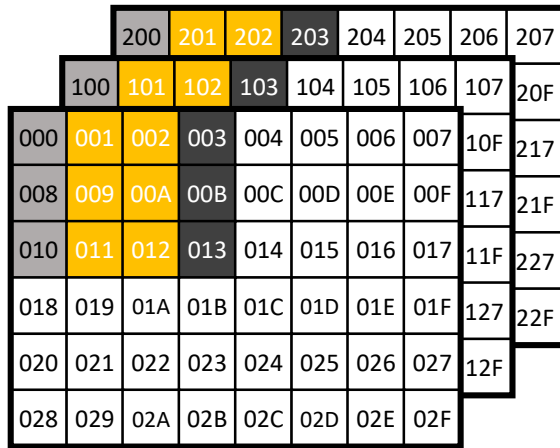
Architectural template



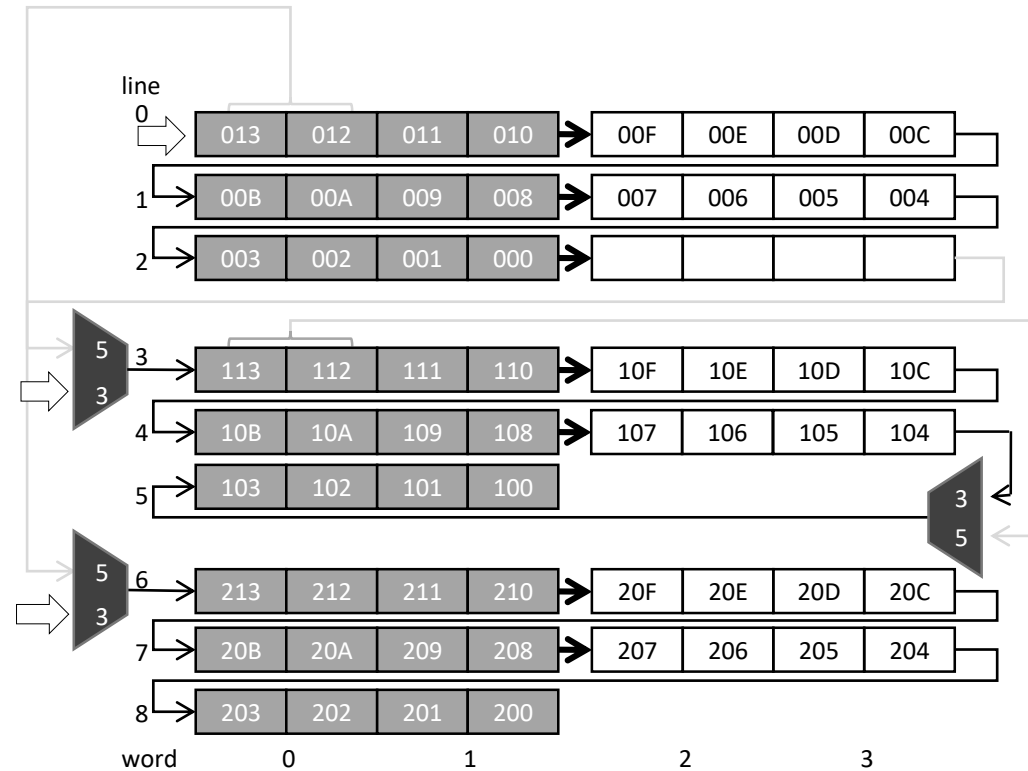
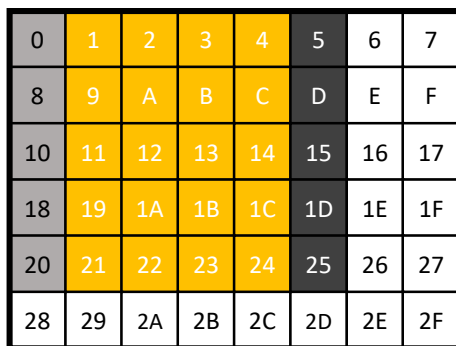
Dynamic reconfiguration support

- Support 5x5 and 3x3 layers
- SoP modules sized to suit both kernel sizes (27 MACs)
 - a) 3x3 filters on 3 input features
 - b) 5x5 filters on 1 input feature (25 MACs)
- Line buffer implements adaptivity

a) 3x3



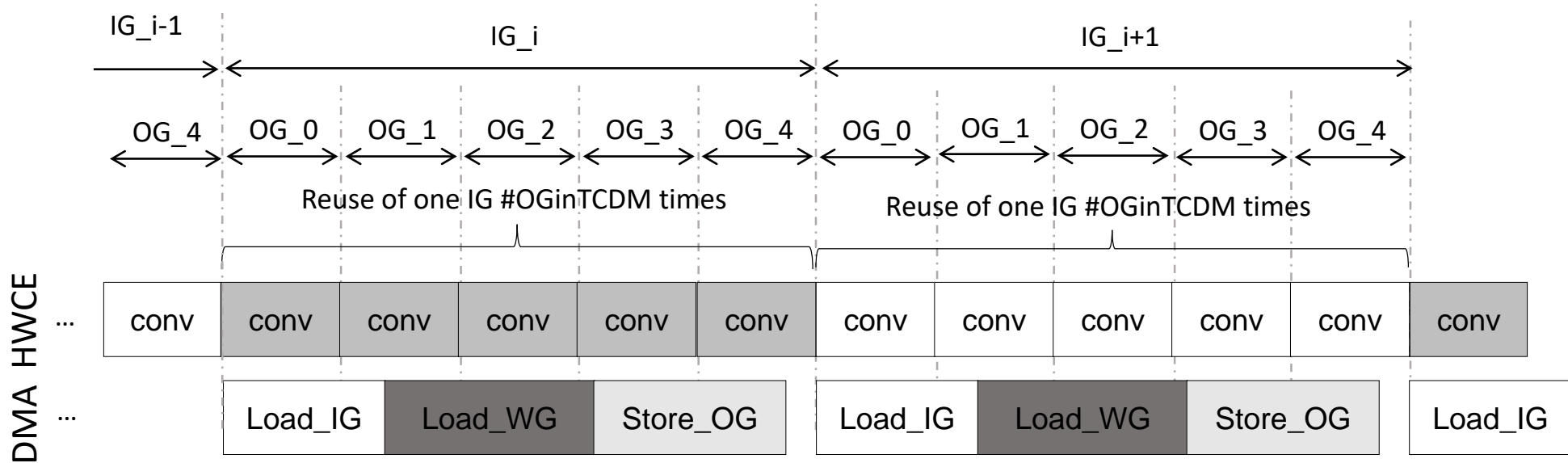
b) 5x5



Convolution scheduling

IG (Input group): 4/12 input features simultaneously loaded by HWCE

OG (Output group): 4 output feature contributions simultaneously produced by HWCE



- Load an IG
- Reuse it! Compute its contributions to several output groups
- Overlap communication and computation (double buffering) – Load next IG and weights during conv
- Iterate until all the contributions to the OGs are accumulated
- Use DMA idle slots to output OGs when complete

Results – Hardware implementation evaluation

Parameter	Value
Slow clock frequency	75 MHz
HWCE clock frequency	150 MHz
Peak performance	129.6 GMAC/s (260 GOPS/s)
I/O bandwidth	16 B/cycle
Line buffer size	128 word
Pixel data precision	16 bit

Resource	DSP	BRAM	LUTs as logic	LUTs as SRegs	Regs
Used	874	192	86047	19989	97585
Available	900	545	218600	218600	437200
Utilization	97.1%	32.2%	39.4%	28.4%	22.3%

Other Works on FPGAs

- Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks, Zhang et al [Zhang 2015]
 - High Level Synthesis approach (loop unrolling and pipelining selected after DSE)
- Going Deeper with Embedded FPGA Platform for Convolutional Neural Network, Qiu et al. [Qiu2016]
 - Model compression (SVD), convolution AND dense layers

	Zhang 2015	Qiu 2016
Platform	Virtex7 VX485t	Zynq XC7Z045
Clock (MHz)	100	150
Quantization	32-bit float	16-bit fixed
Logic Utilization	186K (61%)	183K (84%)
DSP Utilization	2240 (80%)	780 (89%)
BRAM Utilization	1024 (50%)	486 (87%)
Total GOP in Network	1.33	30.8
Performance (GOP/s)	61.6	137.0
Power (W)	18.6	9.6
Energy Efficiency (GOP/J)	3.3	14.2

Binarized Neural Network

- Accelerating Binarized Convolutional Neural Networks with Software-Programmable FPGAs, Zhao et al. [ZhaoFPGA2017]
- FINN: A Framework for Fast, Scalable Binarized Neural Network Inference, Umuroglu et al.

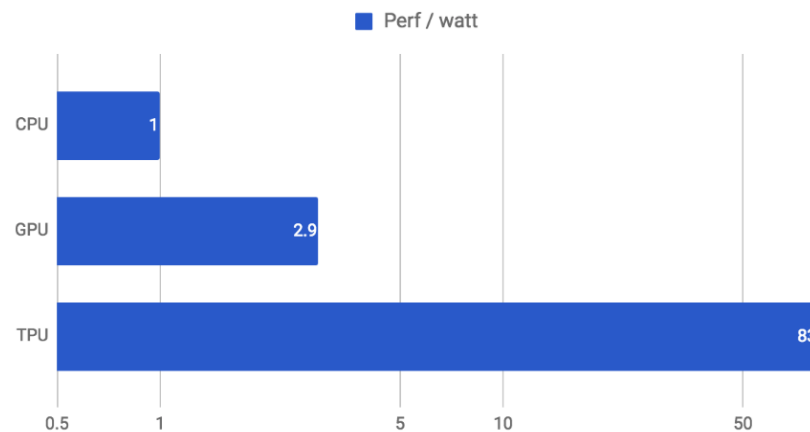
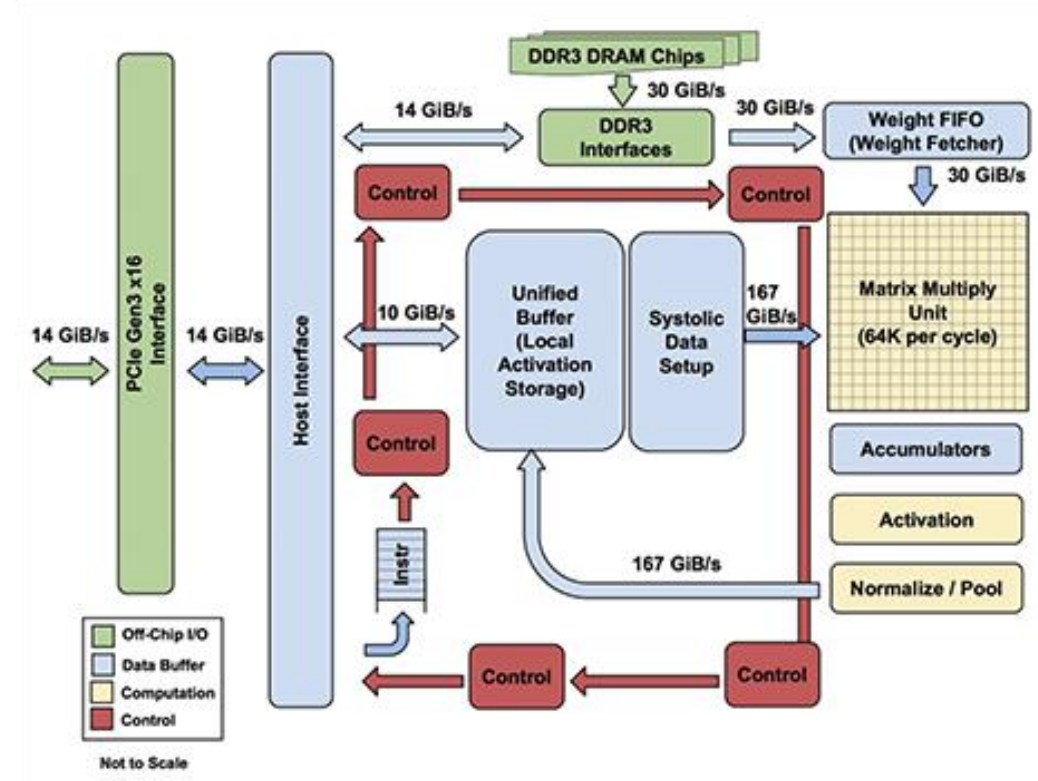
Custom processing platforms

- People who are really serious about software should make their own hardware.

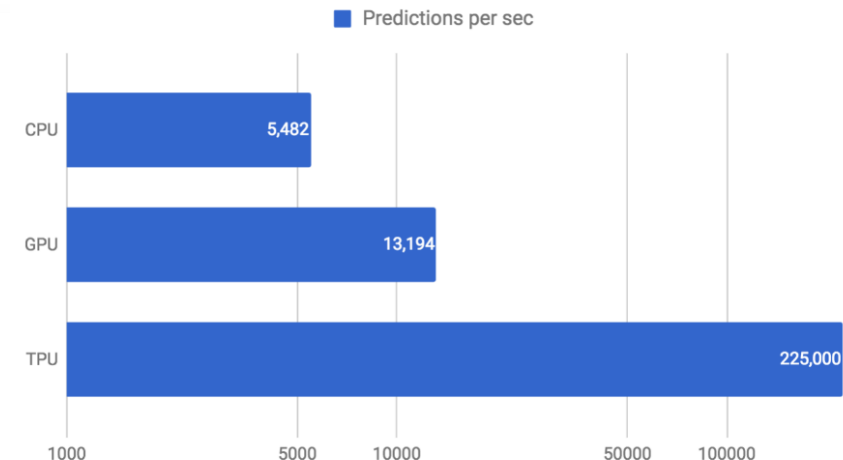
Alan Kay (1982)

Google's Tensor Processing Unit (TPU)

- The TPU includes the following computational resources:
 - **Matrix Multiplier Unit (MXU):** 65,536 8-bit multiply-and-add units for matrix operations
 - **Unified Buffer (UB):** 24MB of SRAM that work as registers
 - **Activation Unit (AU):** Hardwired activation functions
- Controlled by a dozen high-level instructions for neural network inference.



Performance / watt, relative to contemporary CPUs and GPUs (in log scale)(Incremental, weighted mean)



Throughput under 7ms latency limit (in log scale)(99th% response with MLPD: CPU = 7.2 ms, GPU = 6.7 ms, TPU = 7.0 ms)

Other ASIC accelerators/PEs

Publication	Throughput [GOPS]	En.Eff. [GOPS/W]	Supply [V]	Area Effic. [GOPS/MGE]
Neuflow [Pham2012]	320	490	1.0	17
EIE [Moons2016]	102	2600	0.5 - 1.1	64
Eyeriss [Chen2016]	84	160	0.8 - 1.2	46
NINEX [Park2016]	569	1800	1.2	51
k-Brain [Park2015]	411	1930	1.2	109
Origami [Cavigelli2015]	196/74	437/803	1.2/0.8	90/34
YodaNN [Andri2016]	1510/55	9800/61200	1.2/0.6	1135/41

Tools

TensorFlow	Google Brain, 2015 (rewritten DistBelief)
Theano	University of Montréal, 2009
Keras	François Chollet, 2015 (now at Google)
Torch	Facebook AI Research, Twitter, Google DeepMind
Caffe	Berkeley Vision and Learning Center (BVLC), 2013

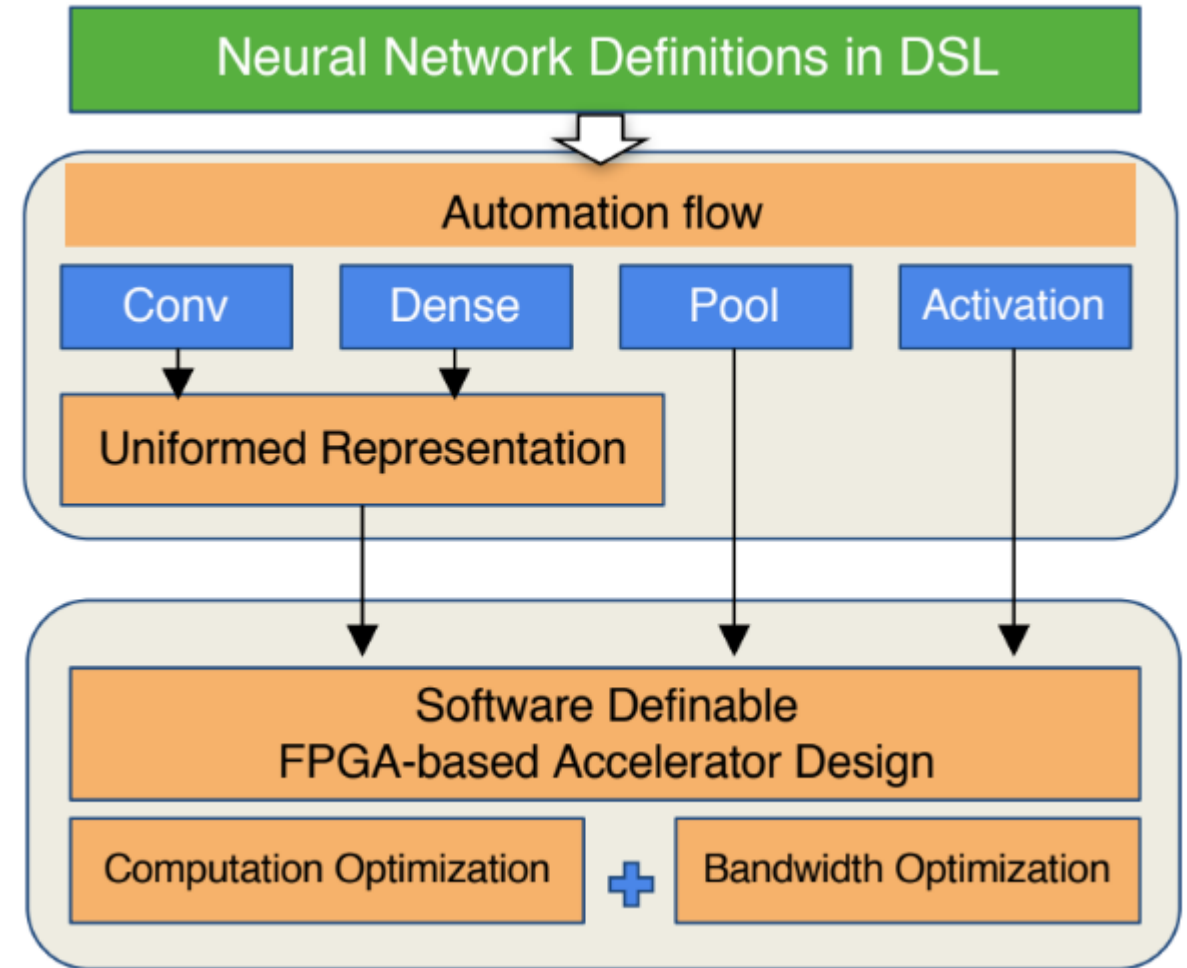
```
layer {
  name: "conv1"
  type: "Convolution"
  bottom: "data"
  top: "conv1"
  param {
    lr_mult: 0
    decay_mult: 0
  }
  convolution_param {
    num_output: 64
    kernel_size: 3
    pad: 1
  }
}

layer {
  name: "loss"
  type: "SoftmaxWithLoss"
  bottom: "fc8"
  bottom: "label"
  top: "loss"
}
```

- Platform-specific Implementation tools
- The NVIDIA CUDA[®] Deep Neural Network library (cuDNN) is a GPU-accelerated library of primitives for [deep neural networks](#).
 - highly tuned implementations for forward and backward convolution, pooling, normalization, and activation layers.
 - part of the [NVIDIA Deep Learning SDK](#).
- Xilinx Deep Neural Network (xfDNN) library is highly optimized for building deep learning inference applications. Designed for maximum compute efficiency at 16-bit and 8-bit integer data types.

Holistic flows

- Caffeine: Towards Uniformed Representation and Acceleration for Deep Convolutional Neural Networks, Zhang et al. ICCAD'16, Nov 2016



Coming soon - ALOHA

