



Challenges and Solutions for Embedding Vision AI

Charles Qi
System Solution Architect
March/25/2018
V 0.4

cadence®

Deep Learning is Broad

Wide spectrum of applications, wider spectrum of technology

Applications

Social media and big-data search:

Google, Facebook, Microsoft, Baidu, NEC, IBM, Yahoo, AT&T



Medical: genomics, radiology, screening, protein sequencing



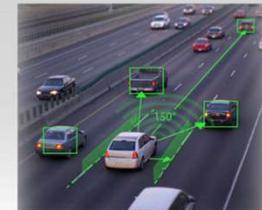
Finance: TradeTrek, M.J. Futures, Alyuda



Speech recognition: Apple, Google, Nuance, Microsoft



Vision/ADAS: Nvidia, Intel



Public Frameworks

Caffe



theano

mxnet

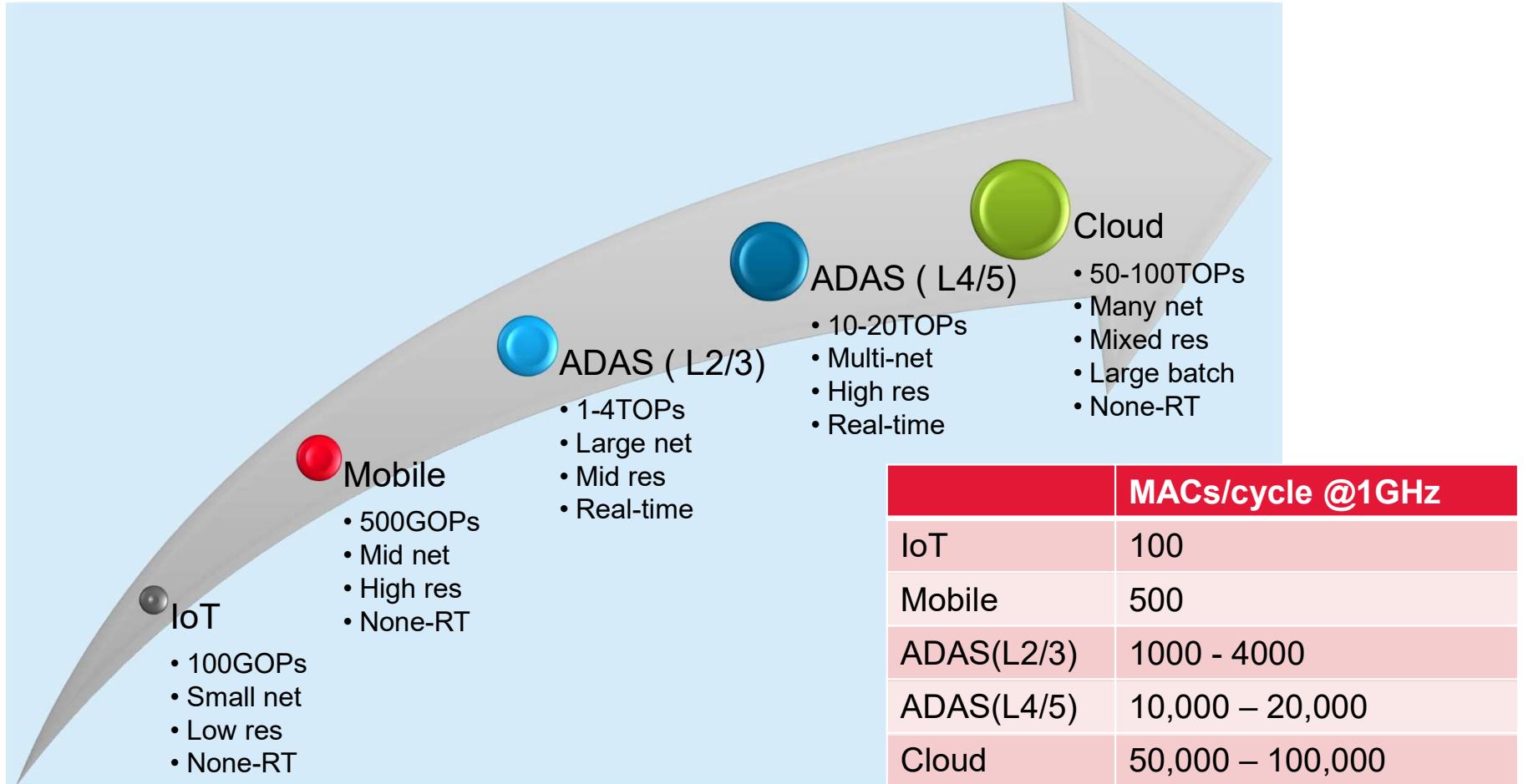


PYTORCH

HW Platforms

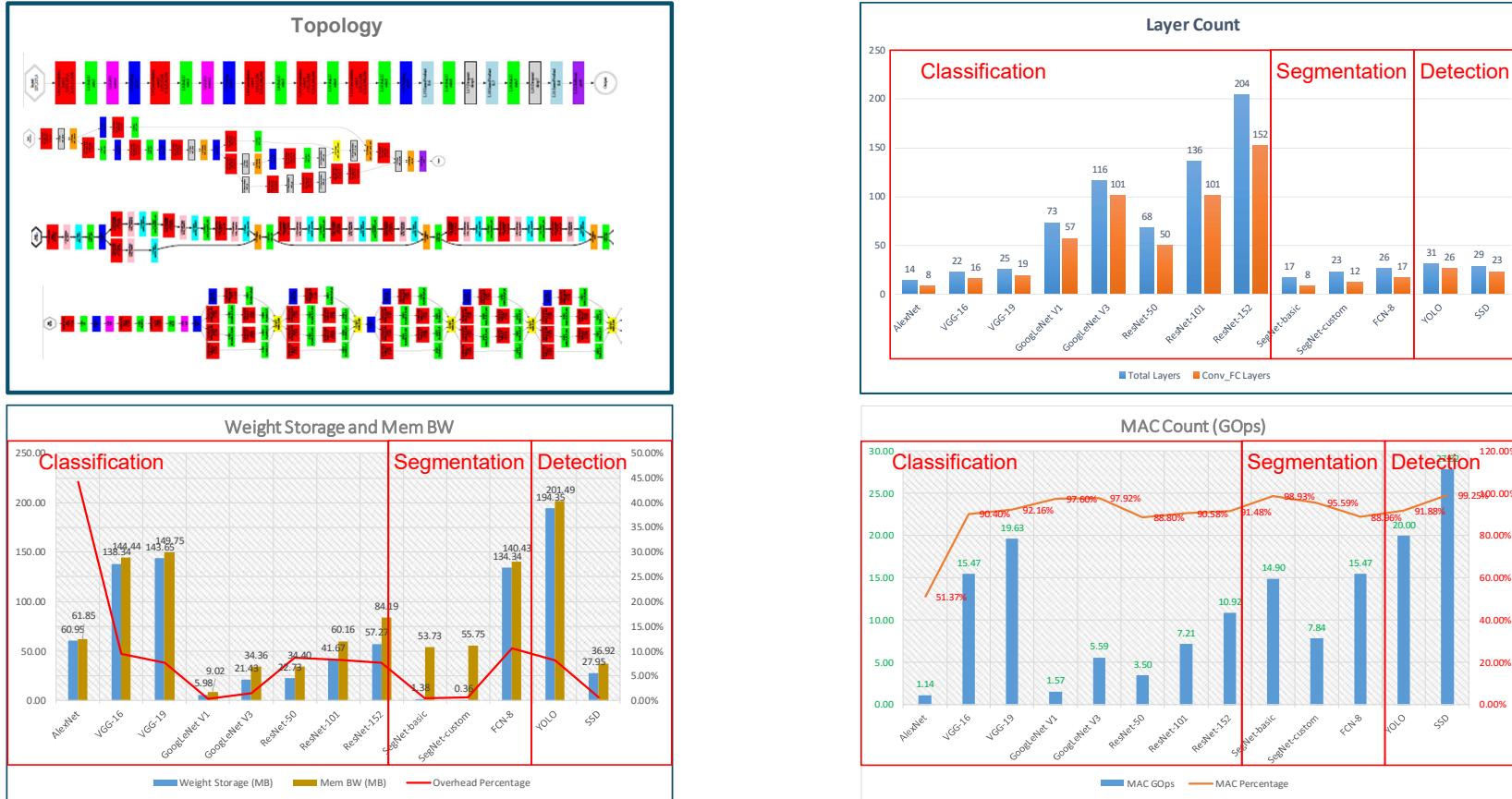


CNN: Convolution Performance is Key – How to Scale?



Diversity Drives Neural Network Complexity

Net topology, layer count, feature map and kernel sizes differ significantly



Conv and FC layers account for 90% load

Convolution Operation Deep Dive

- Convolution layer is multi-dimensional
 - a 3-D input feature tensor convolving with a set of 3-D filter tensors
 - Deep loop stack in software
 - Large variation of input spatial dimension, channel depth and filter count

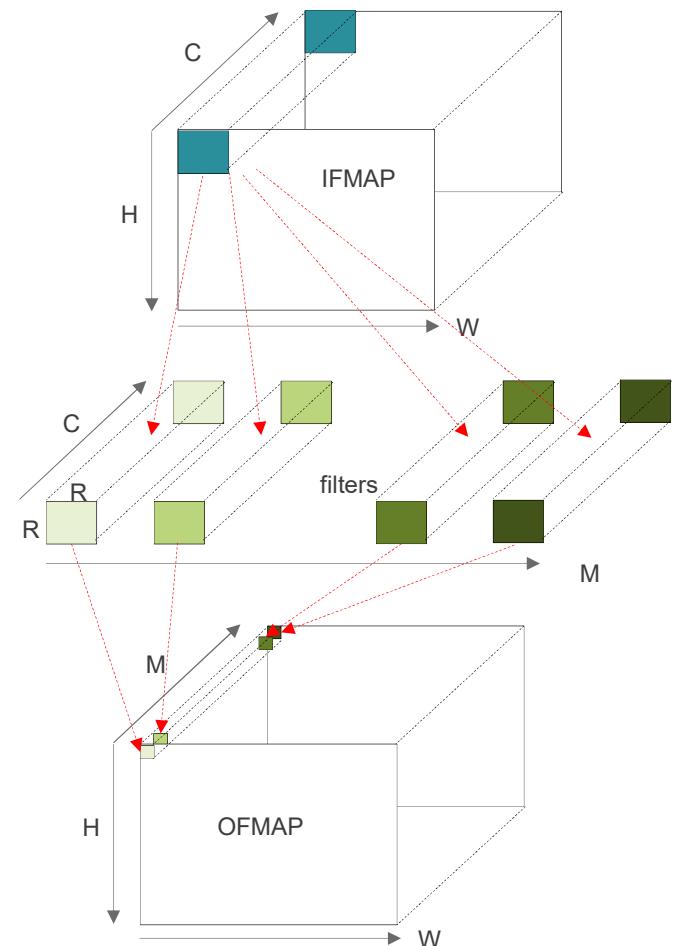
Operation:

$$F_o(x, y, n) = \sum_{c=0}^{C-1} \sum_{r=-R}^R \sum_{r=-R}^R W(r, r, c) * F_l(x + r, y + r, c) \quad \forall x \in X, y \in Y, m \in M$$

Batch → Implementation as 8 nested loops:
 Layers → `for (b=0; b<B; b++)`
 Filter set → `for (l=0; l<L; l++)`
 Ifmap height → `for (m=0; m<M; m++)`
 Ifmap width → `for (y=0; y<H; y++)`
 Channels → `for (x=0; x<W; x++)`
 Kernel width → `F_o(x, y, m) = 0`
 Kernel height → `for (c=0; c<C; c++)`

`for (ry=0; ry<RY; ry++)`
`for (rx=0; rx<RX; rx++)`

$$F_o(l, b, x, y, m) += W_{l, b, m}(rx, ry, c) * F_l(x + rx, y + ry, c)$$



GEMM Approach Non-ideal for Embedded AI

- GEMM converts convolution into 2-D matrix multiply by brute-force
- IFMAP
 - TxN array, linearize spatial dimension X*Y rows
 - Replication of activation data by R*R across C channels
- Filter
 - N*M array, R*R*C channel element rows, M filter columns

Input Feature Map Matrix

	$N = R \times R \times C$											
	(x0,y0)	(x1,y0)	(x2,y0)	(x0,y1)	(x1,y1)	(x2,y1)	(x0,y2)	(x1,y2)	(x2,y2)	(x3,y2)	(x4,y2)	
t=0	c0	c1	c0	c1	c1	c0	c1	c0	c1	c1	c0	c1
t=1	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1
t=2	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1
t=3	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1

t=T-4	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1
t=T-3	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1
t=T-2	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1
t=T-1	c0	c1	c1	c0	c1	c1	c0	c1	c1	c0	c1	c1

$T = X \times Y$

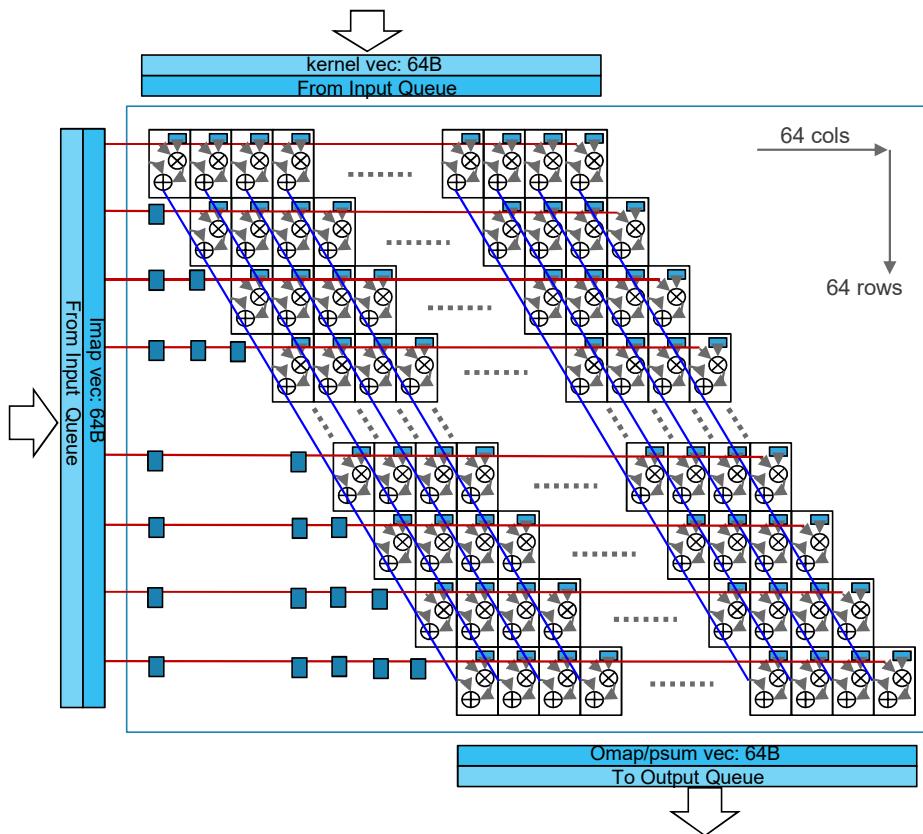
Filter Matrix

	M											
	m	m	m	m	m	m	m	m	m	m	m	m
	=	=	=	=	=	=	=	=	=	=	=	=
	0	1	2	3	-	-	-	-	-	-	-	-
lx=0	c0	c0	c0	c0	c1							
ky=0	c1											
lx=1	c0	c0	c0	c0	c1							
ky=0	c1											
lx=2	c0	c0	c0	c0	c1							
ky=0	c1											
lx=1	c0	c0	c0	c0	c1							
ky=1	c1											
lx=0	c0	c0	c0	c0	c1							
ky=2	c1											
lx=1	c0	c0	c0	c0	c1							
ky=1	c1											
lx=2	c0	c0	c0	c0	c1							
ky=0	c1											
lx=1	c0	c0	c0	c0	c1							
ky=2	c1											
lx=2	c0	c0	c0	c0	c1							
ky=1	c1											
lx=0	c0	c0	c0	c0	c1							
ky=2	c1											
lx=1	c0	c0	c0	c0	c1							
ky=2	c1											

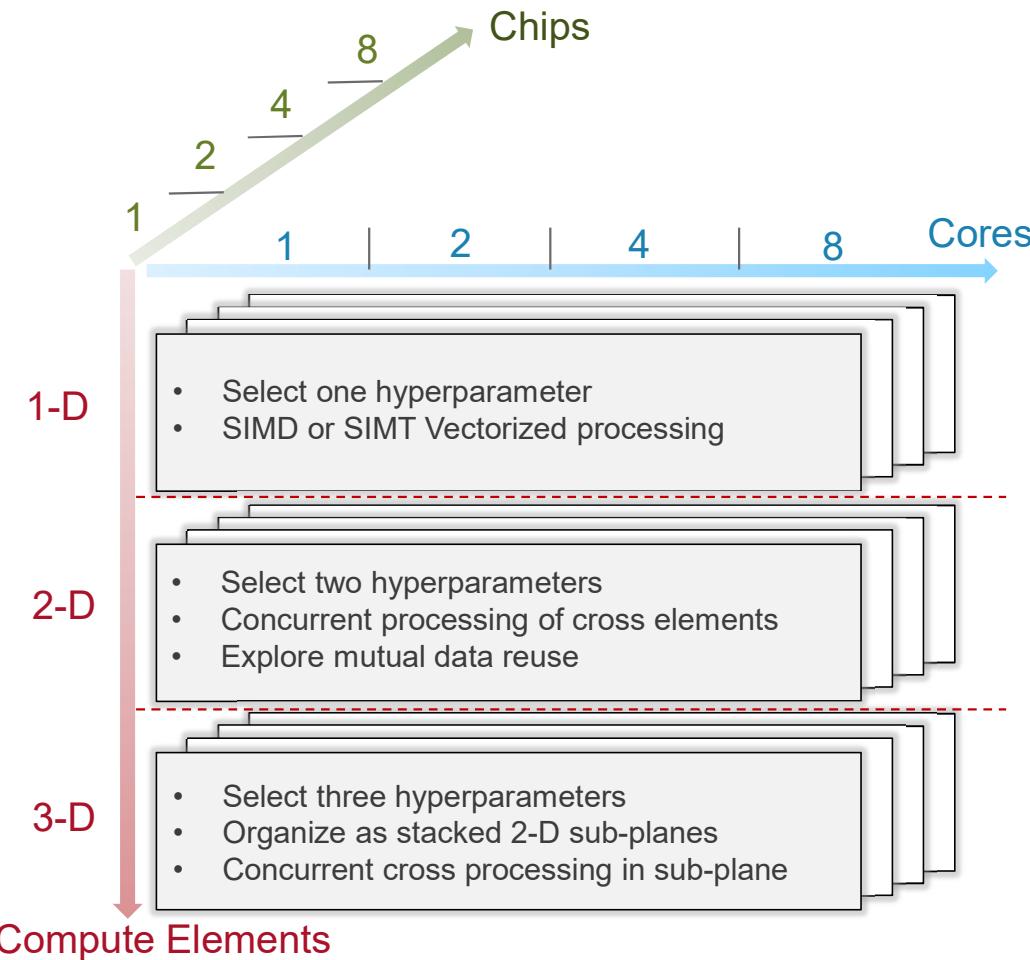
$N = R \times R \times C$

NxN Plain GEMM Accelerator

- NxN, N=64 MAC array
 - 64 independent vertical slices
 - 64 accumulator rows
- Stationary weight array
 - 64 rows of kernel weights preloaded to array cells in 64 cycles
- Dynamic input/output vectors
 - One 64B/cyc imap vector is fed to array
 - One 64B/cyc omap vector is produced
- Psum stationary array also possible
- Limitations
 - Rigid array configuration often underutilized in various layer dimensions
 - High pipeline overhead
 - im2col implodes memory footprint
 - Inefficient for vector*vector, vector*matrix without large batch
 - Hard to partition for multi-core



Scaling the Dimensions to Accelerate



Nominal network size:

- 50 layers
 - 500M MACs
 - 60 fps
- 1.5T MACs

1D vector engine (GPU, DSP or HW):

- Frequency: 1GHz
- 1000 – 2000 MACs/cycle
- At the limit of datapath width
- Utilization hard to sustain

2D or 3D array engine:

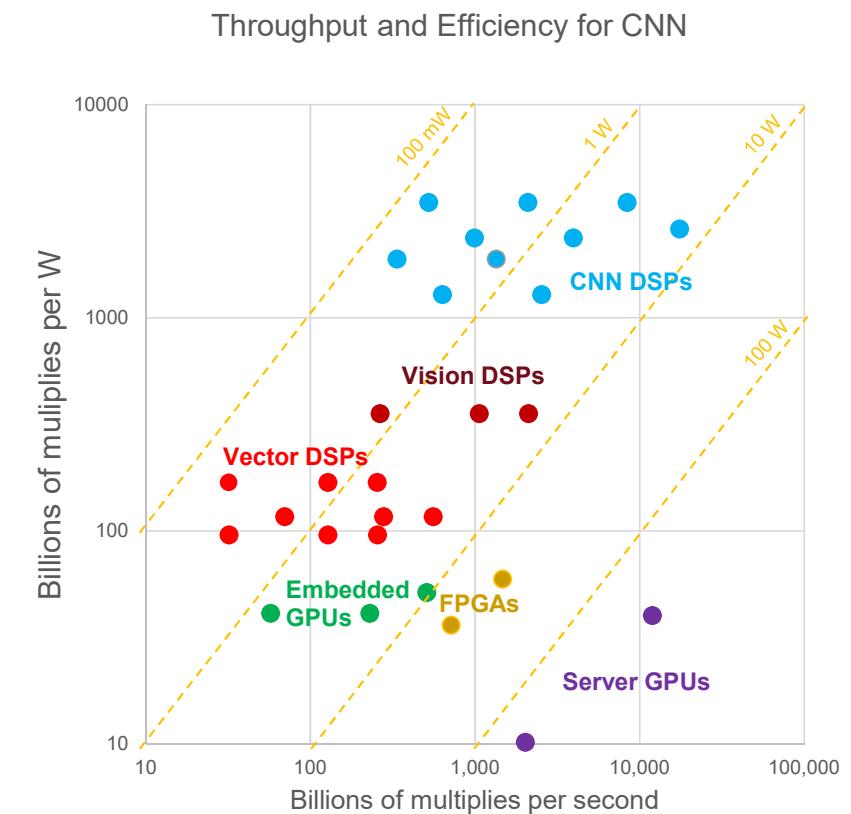
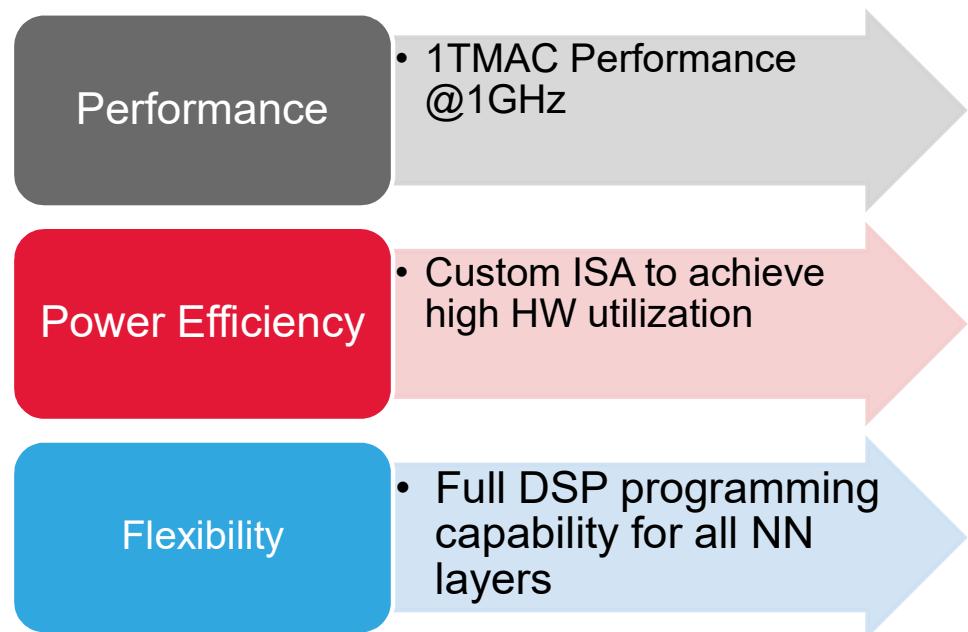
- Frequency: 0.5-1GHz
- 4K MACs/cycle
- 2D: 64x64
- 3D: 32x32x8
- Easy to achieve high utilization
- Data reuse minimize IO
- More complex scheduling

Extract Parallelism for Multi-dimensional Scaling

Dimension	Dep.	Single Core Multi-Dimension Vectorization	Multi-core (U) Parallel Load Split
Batch of images (B)	Indep.	<ul style="list-style-type: none"> Typically not applied because $B \ll V$ May be applicable to narrow SIMD Requires interleaving data from multiple “images” 	<ul style="list-style-type: none"> Easy to distribute smaller batch to cores when $B \gg U$ Reduce filter weight reuse (weight replication, reduced W reuse) Latency to process one image not reduced, bad for RT application
Layer (L)	Data dep	<ul style="list-style-type: none"> NA, layers processed sequentially or pipelined in different cores due to layer-to-layer data dependency 	<ul style="list-style-type: none"> No overlapping weights in cores, most effective weight reuse Challenge is to balance compute load for different layers Core-to-core data dependency requires core-to-core pipelining
Number of kernels (M)	Indep.	<ul style="list-style-type: none"> Most efficient for wide SIMD ($M \gg V$ and $M \% V == 0$) Allow kernel weight or data reuse Large local mem footprint due to parallel accumulations Requires data and weight reorganization 	<ul style="list-style-type: none"> Easy to partition for U cores without weight overlap when $M \gg U$ May limit core-level SIMD vectorization on M when $M < V * U$ Input data is replicated across cores
Spatial (X,Y)	Indep.	<ul style="list-style-type: none"> Typically applied to X only when $X \gg V$ Filter weight can be reused for entire vector Inefficiency if $X \% V$ is non-zero 	<ul style="list-style-type: none"> Flexible to partition tiles to cores Inefficiency when tile size becomes smaller than SIMD width Filter weights replicated in cores processing different tiles
Input depth (C)	Accum.	<ul style="list-style-type: none"> Inefficient for shallow layer with $C \ll V$ Require vector reduction of the partial accumulator values in SIMD ways 	<ul style="list-style-type: none"> Ineffective for shallow layer with $C \ll U$ Require final accumulation of intermediate results across cores
Kernel (R)	Accum.	<ul style="list-style-type: none"> Typically not applied due to odd and small R Different rows of R in the same vector computes different data, require data re-org, lack of reuse 	<ul style="list-style-type: none"> Typically not applied due to odd and small R Require final accumulation for R values from different cores

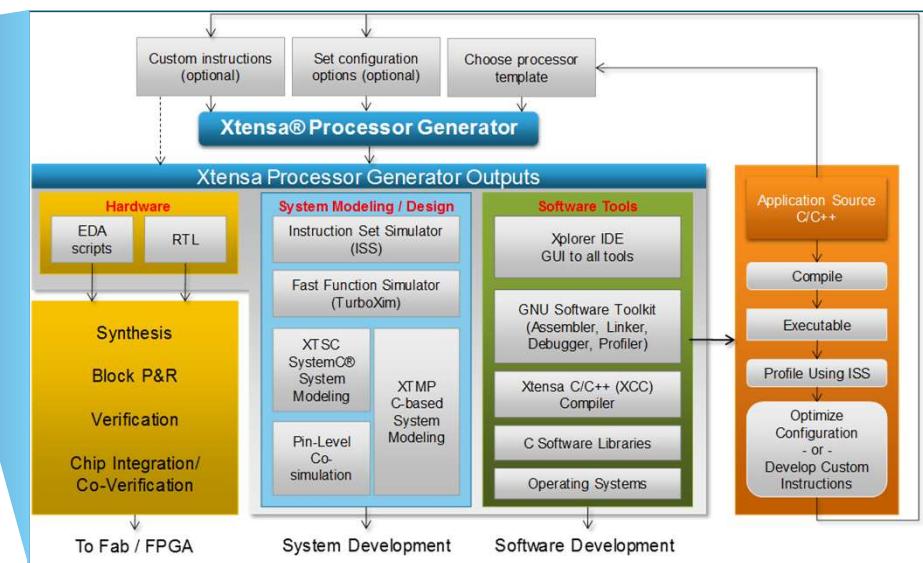
Accelerating Embedded NN at Edge

solution need to achieve high performance with architecture flexibility



Viable Solution: Tensilica Scalable DSP Platform

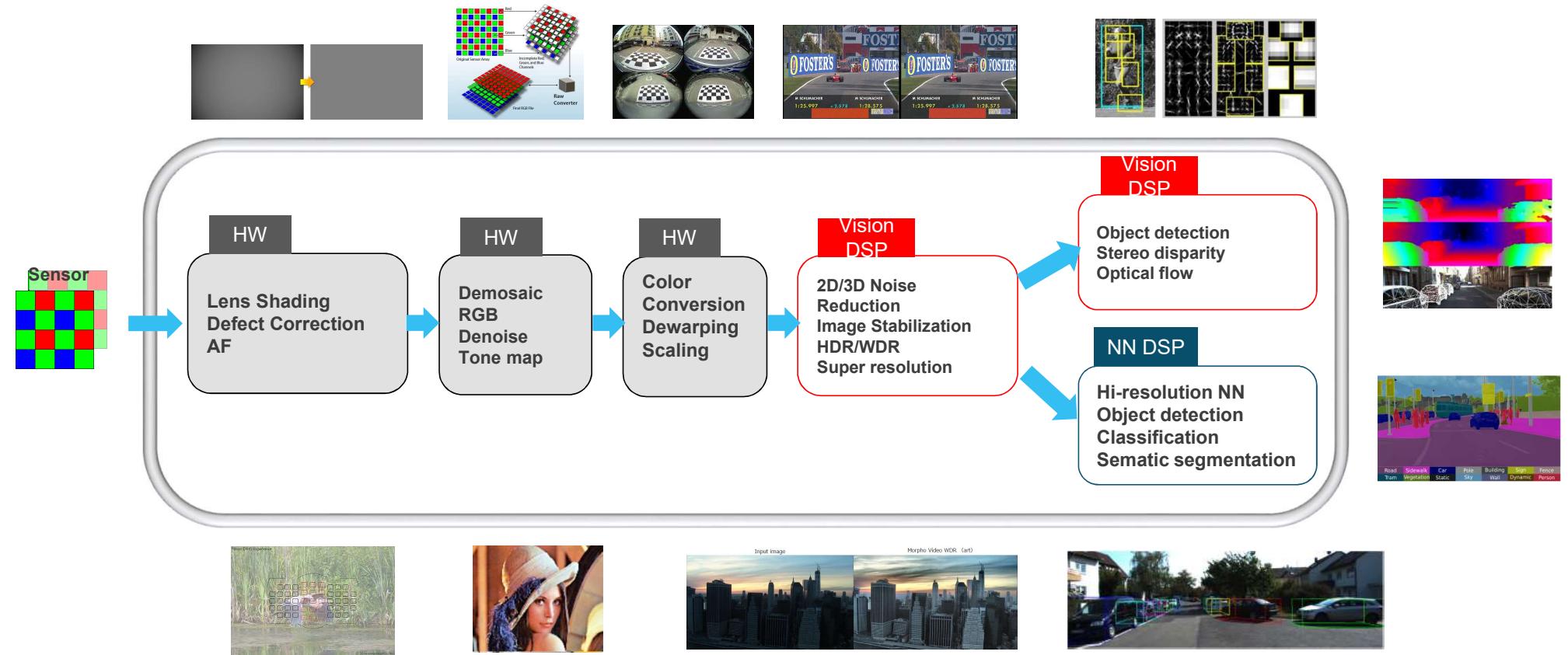
DSP cores tailored for application specific needs



- Scalable performance, consistent programming model
- Common baseline Xtensa architecture + application specific processing optimization
- Fully integrated HW/SW tool chain from user-friendly Xplorer IDE

The Real-time Embedded Vision AI Pipeline

where intelligence and programmability meet



How is DSP Optimized for Embedded Vision/CNN?

Focus on computation efficiency and minimize memory access

Increase computation efficiency

- ✓ ISA level parallelism based on fixed point SIMD
- ✓ Balance memory and compute resource with VLIW
- ✓ Specialized TIE instructions for special tasks

Reduce system memory access overhead

- ✓ Improve access efficiency with local RAM
- ✓ Hide access latency with tiling and DMA
- ✓ Minimize cycle count for non-contiguous access

Example of Special Instruction Extension (TIE)

Achieve significant performance while reducing power

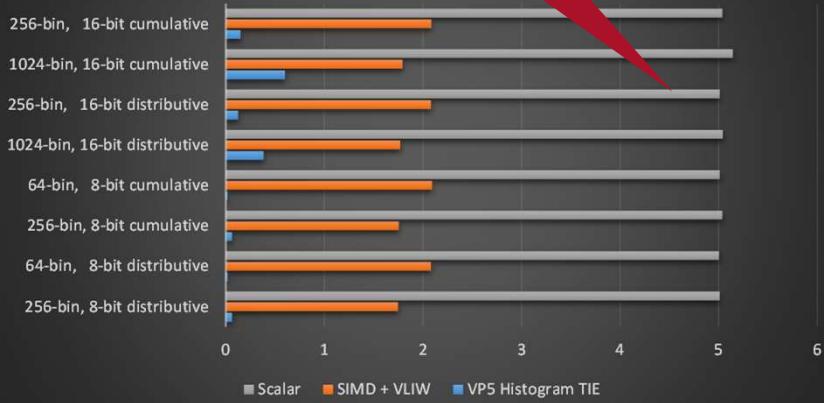
Scalar C

```
//scalar C code
for (idxData = 0; idxData < len; idxData += 4) {
    data = (uint8_t) datain[idxData];
    hist[data >> ShiftRight]++;
    data = (uint8_t) datain[idxData + 1];
    hist[data >> ShiftRight]++;
    data = (uint8_t) datain[idxData + 2];
    hist[data >> ShiftRight]++;
    data = (uint8_t) datain[idxData + 3];
    hist[data >> ShiftRight]++;
}
```

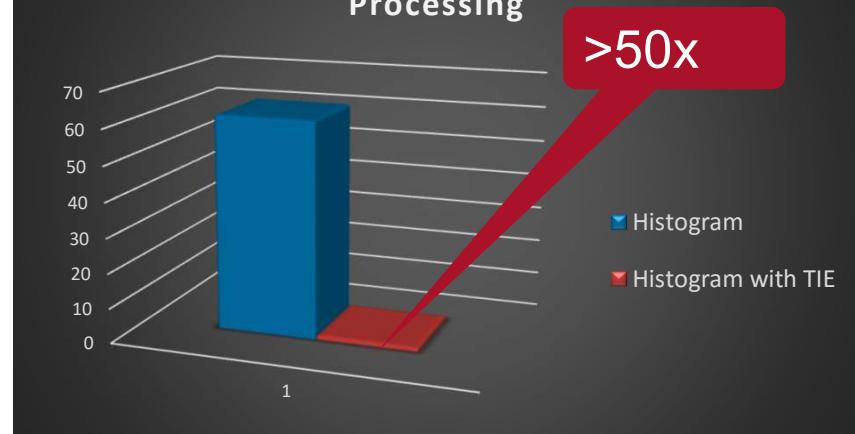
Specialized TIE for histogram

```
//specialized TIE created for histogram processing
for (int y = 0; y < height; ++y) {
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists0 = IVP_ADDNX16(vhist0, vdst); // 0 .. 31
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists1 = IVP_ADDNX16(vhist1, vdst); // 32 .. 63
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists2 = IVP_ADDNX16(vhist2, vdst); // 64 .. 95
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists3 = IVP_ADDNX16(vhist3, vdst); // 96 .. 127
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists4 = IVP_ADDNX16(vhist4, vdst); // 128 .. 159
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists5 = IVP_ADDNX16(vhist5, vdst); // 160 .. 191
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists6 = IVP_ADDNX16(vhist6, vdst); // 192 .. 223
    IVP_COUNTEQ4NX8(vdst, sr, vec0, vec1); vhists7 = IVP_ADDNX16(vhist7, vdst); // 224 .. 255
}
```

Scalar vs. SIMD+VLIW vs. Special TIE
Performance



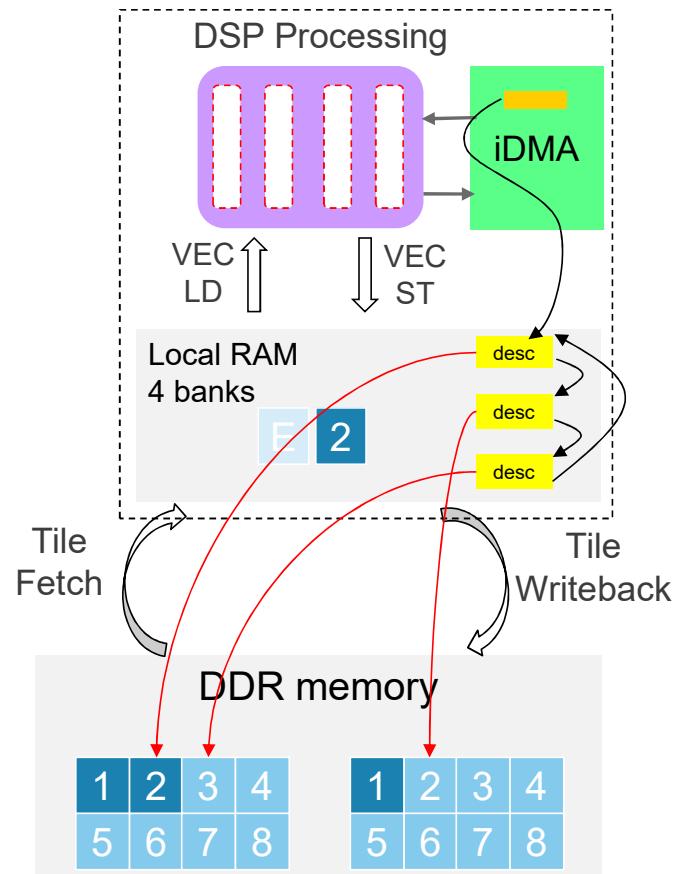
Energy in UJ for VGA Resolution Processing



Hiding Memory Access Latency via Tiling and DMA

Offload tile fetch and writeback

- Tile fetch and writeback by software is inefficient, blocking data processing
- Build-in DMA performs data moving in parallel
 - Parallel DMA access to local memory banks
 - Dedicated NoC master port for DMA
 - Multiple outstanding requests in pipe to memory
- 2D descriptor ring supports complex memory access patterns
 - DMA is aware of the row-dominant image storage format by skipping row pitch at the end of each row
 - 3D/4D tile stream can be composed using multiple descriptors
 - Flexible data alignment
- DMA register-mapped inside processor to minimize initiation/termination overhead
- Ping-pong tile buffers in local RAM removes dependency



Representative NN DSP Architecture Considerations

Uniquely design ISA to strive for processing efficiency

256 – 1024 8x8 multiple-accumulate operations per cycle

Vector-by-vector, vector-by-scalar operations with multiple accumulators

Flexible vectorization scheme to maximize HW utilization

Special addressing for load/store 3D tensor data

Acceleration for pooling, nonlinear activation layers

On-the-fly Coefficient Compression/Decompression

NN ISA Architecture Optimization Concept (1)

Amortize IO cost across multiple computations

- Problem: Vector x Vector operations requires 2 SIMD loads + 1 SIMD store
- Explore Data Reuse prosperity in CNN
 - Multiple filters applied to common IFMAP input data
 - Multiple spatial locations of IFMAP convolve with same filter elements
- Optimization: perform vector x scalar for multiple scalars
 - Amortize vector load BW and keep computation pipeline full

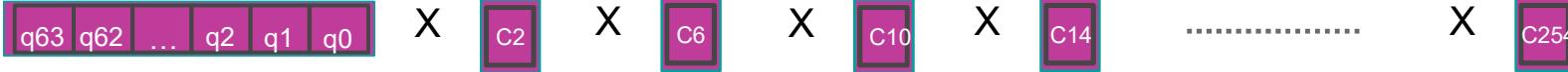
Vec load 0



Vec load 1



Vec load 2



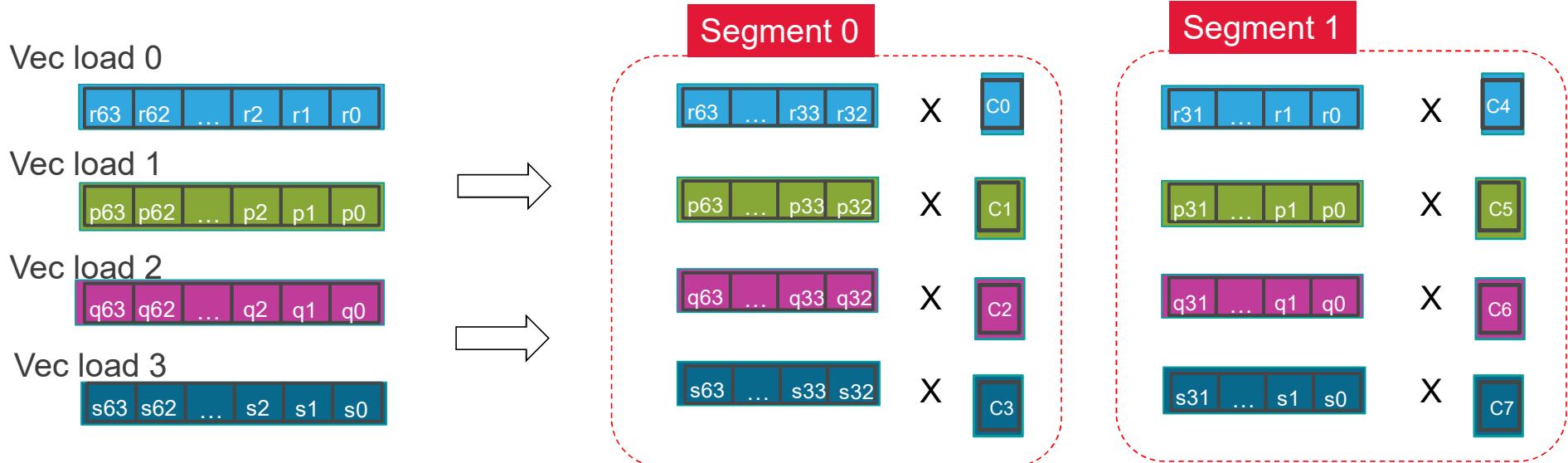
Vec load 3



NN ISA Architecture Optimization Concept (2)

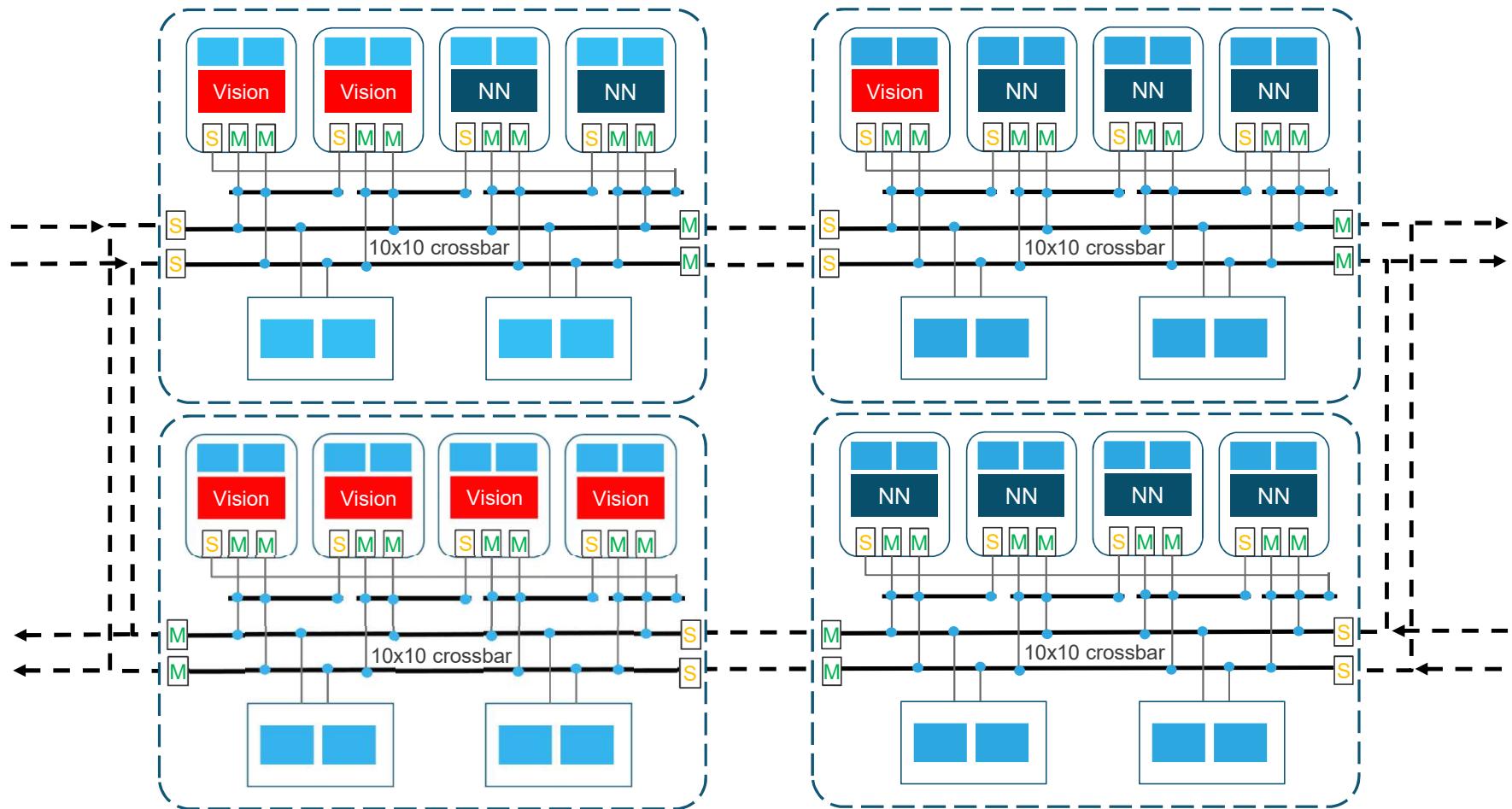
SIMD segmentation

- Problem: layer dimension much smaller than SIMD width
- Optimization: break wide SIMD into smaller segments
 - Apply different scalar data to different segments to fully utilize HW resource



Scale Vision Sub-system Heterogeneous Multi-core

Flexibility to customize MP cluster under the same programming model



Multi-core NN Load Partition Example

Split load across layer/batch/kernel

Operation:

$$F_o(x, y, n) = \sum_{z=0}^{Z-1} \sum_{r=-R}^R \sum_{r=-R}^R W(r, r, z) * FI(x + r, y + r, z) \quad \forall x \in X, y \in Y, n \in N$$

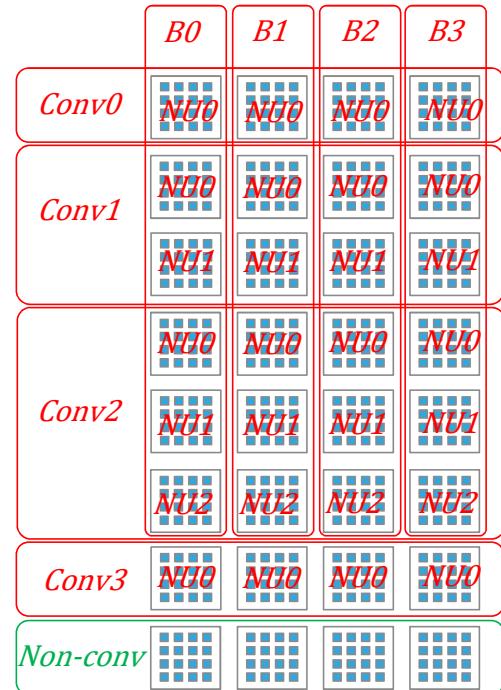
Implementation as nested loops:

```

layers      → for (l=0; l<L; l++)
Ifmap batch → if(layers[l] == CONV) {
    for (b=0; b<B; b++)
kernels     → for (nu=0; nu<NU; nu++)
Ifmap height → for (nv=0; nv<NV; nv++) ← SIMD vectorization in NV or X
Ifmap width  → for (y=0; y<Y; y+=S) ←
Ifmap channels → for (x=0; x<X; x+=S) ←
Kernel width → for (z=0; z<Z; z++)
Kernel height → for (ky=0; ky<KY; ky++)
                for (kx=0; kx<KX; kx++)
                 $F_o(x, y, n) += Wl_{b, n}(kx, ky, z) * FI(x + kx, y + ky, z)$ 
} #end if

```

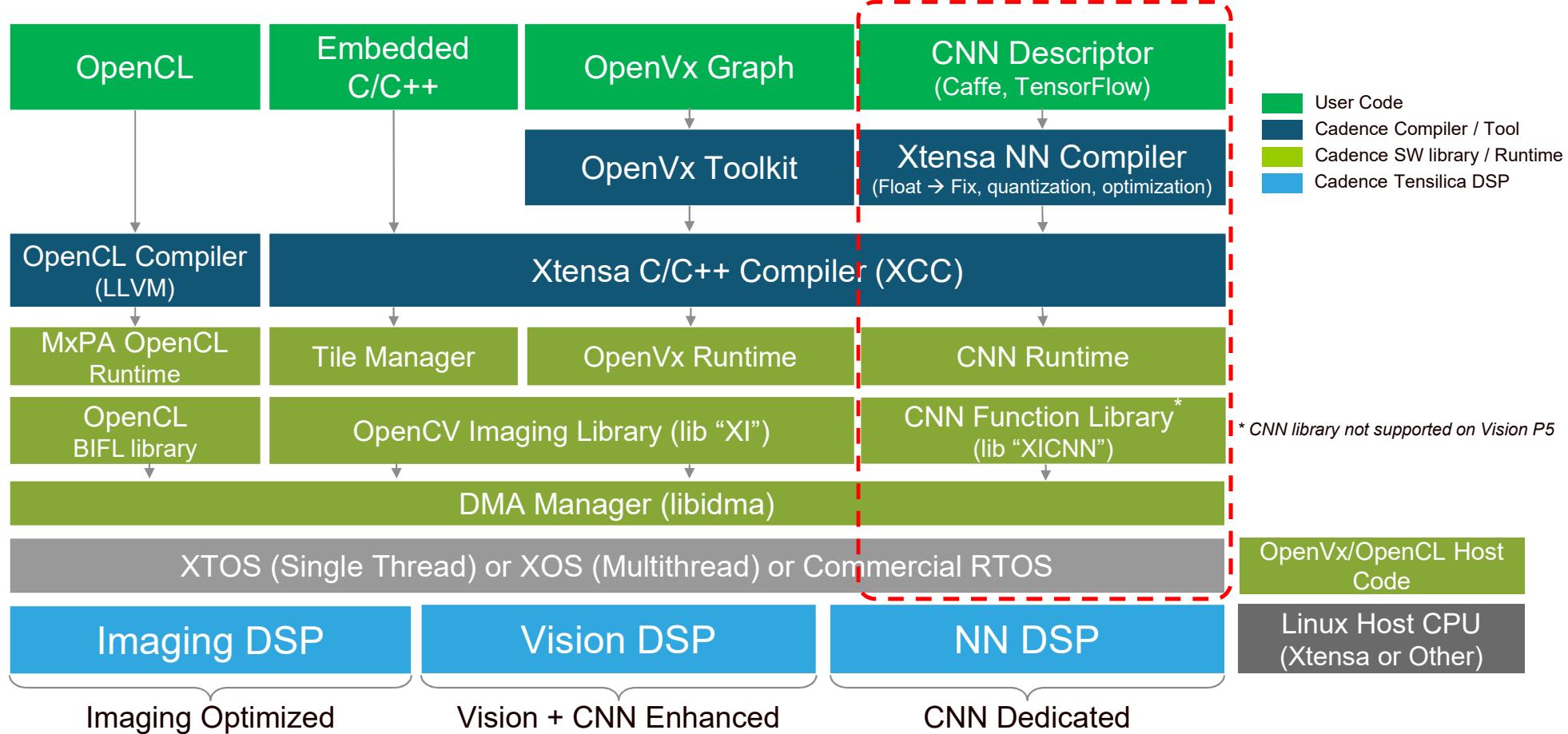
- L and B loops are distributed across MP cores
- N is split into two loops, NU, NV and $N = NU * NV$
 - NU is distributed across cores
 - NV is handled by vectored SIMD



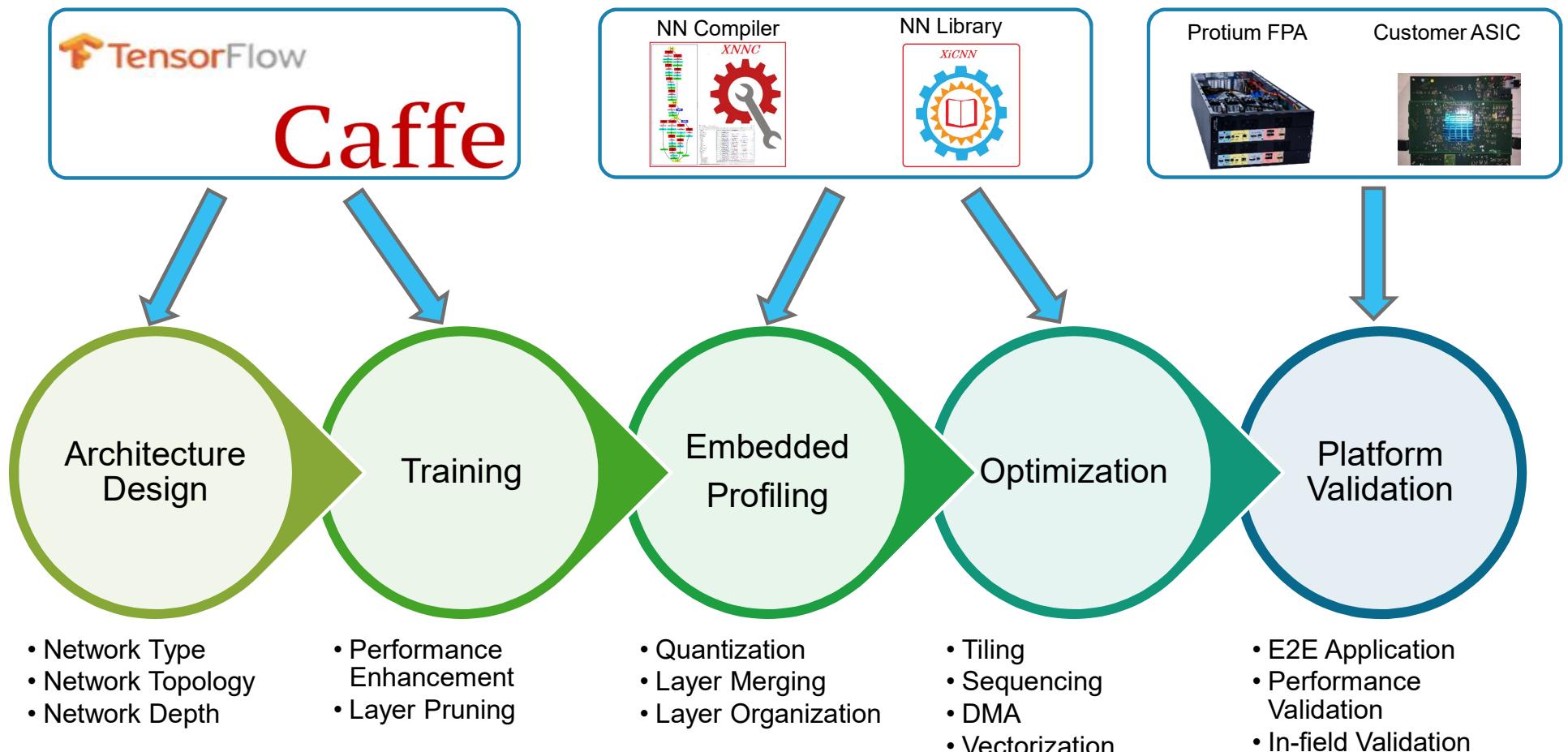
Cores are designated to certain layers and batches. Most efficient if layers can be load-balanced, no overlap in weight coefficients across cores.

Consistent and Comprehensive Vision Software Solution

Full ecosystem of software frameworks and toolchain across common architecture



Concept to Hardware - Embed AI in Five Streamlined Steps



Conclusion



Wide demand for NN in embedded applications

- Diversified application needs drive various NN architectures
- Demand for higher and higher computation performance
- Embedded power profile and implementation cost drive optimization

Optimization of embedded vision/AI remains challenging

- Too many architecture variants makes it hard to optimize across
- Brute-force approach leads to higher power/memory BW, lower utilization
- Need multi-dimensional scaling without losing flexibility

Specialized NN DSP provides efficiency and flexibility

- Highly optimized ISA improves NN computation efficiency
- End-to-end full network implementation in image/vision pipeline
- SW framework and library automates NN implementation

cādēnсē[®]

© 2018 Cadence Design Systems, Inc. All rights reserved worldwide. Cadence, the Cadence logo, and the other Cadence marks found at www.cadence.com/go/trademarks are trademarks or registered trademarks of Cadence Design Systems, Inc. All other trademarks are the property of their respective holders.