

Xilinx ML-Suite

Rahul Nimaiyar
Ashish Sirasao



FPGA for Deep Learning



Deep Learning on Xilinx Adaptable Devices

Data Parallel

- 2D Array of MACs
- Flexible on-chip memory access
- High Bandwidth, Multiple Access Ports

Custom Memory Hierarchy

- Near Memory Compute
- Programmable routing for data & filter reuse

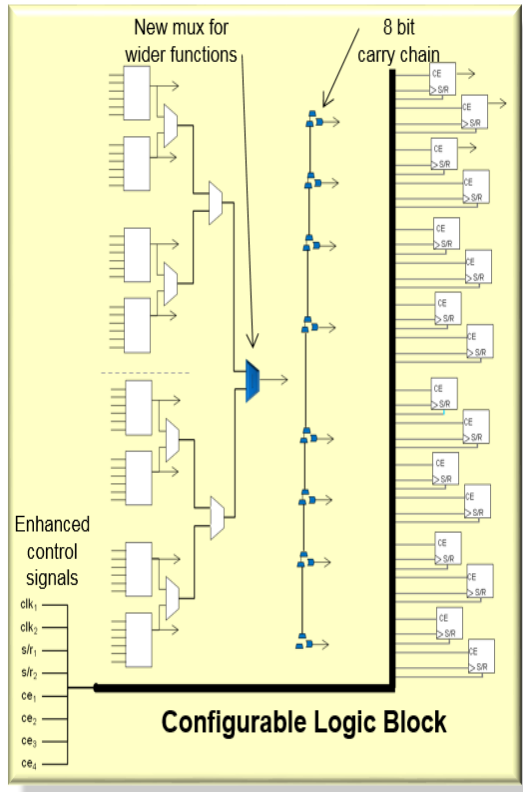
Compression & Sparsity

- Flexible Data Types
- FP32/16, INT16/8/4/2, Binary/Ternary
- Sparsity friendly compute

Broad Device Range

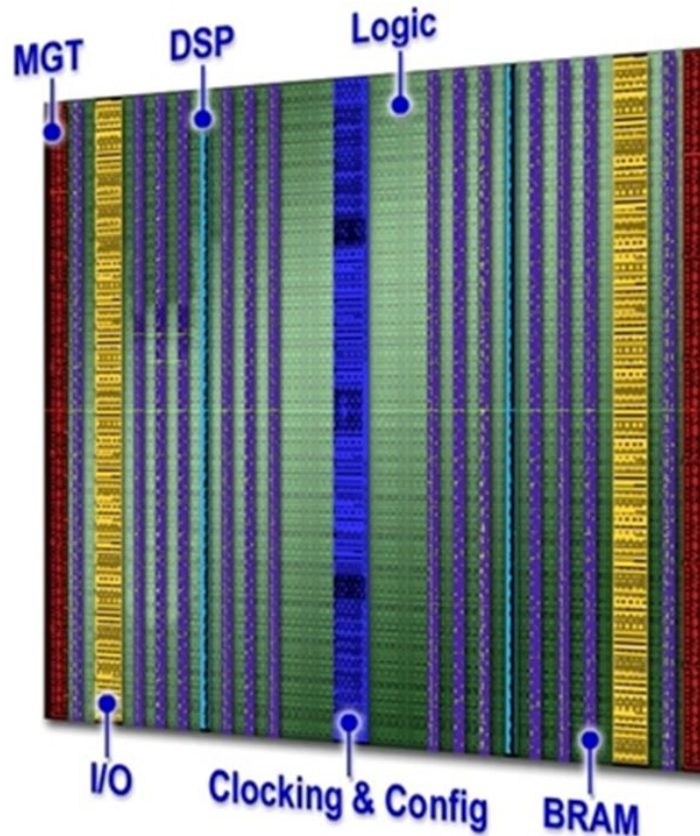
- Scalable device family for different applications
- Built in System functions – Networking, Video, ARM

Xilinx FPGA Building Blocks

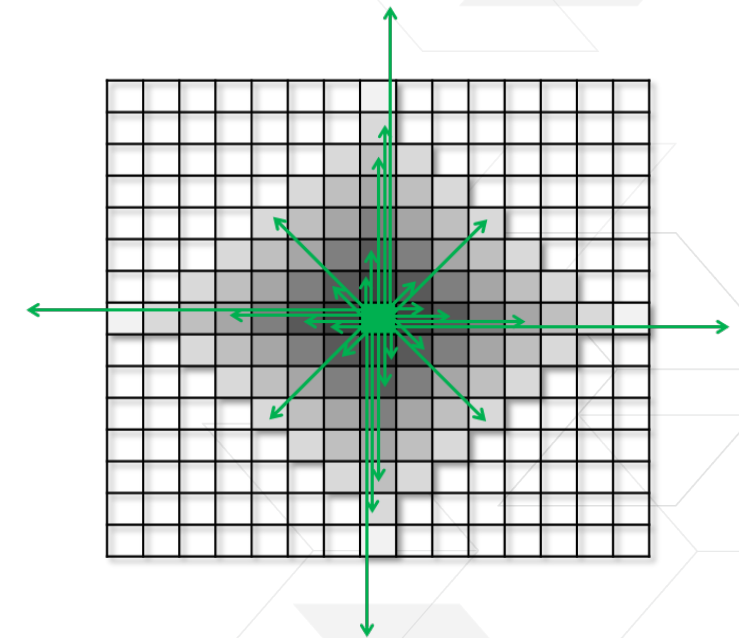


CLB Architecture

Eight six-input LUTs, Sixteen Flip Flops,
Carry chain, Muxes for larger logic
structures, Distributed memory, Shift
register



Columnar Chip Layout

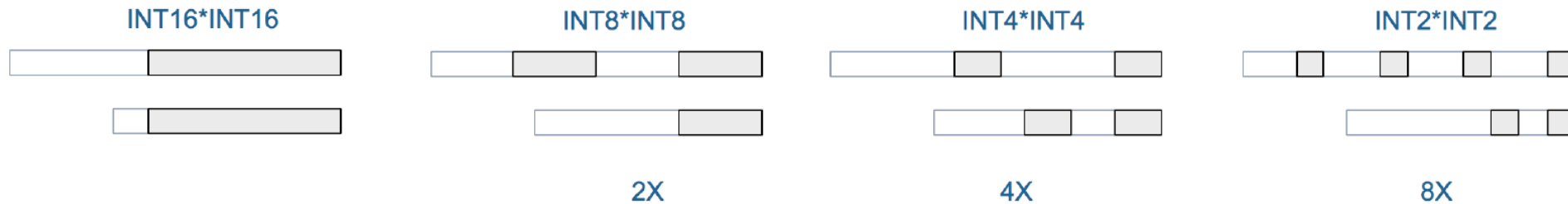
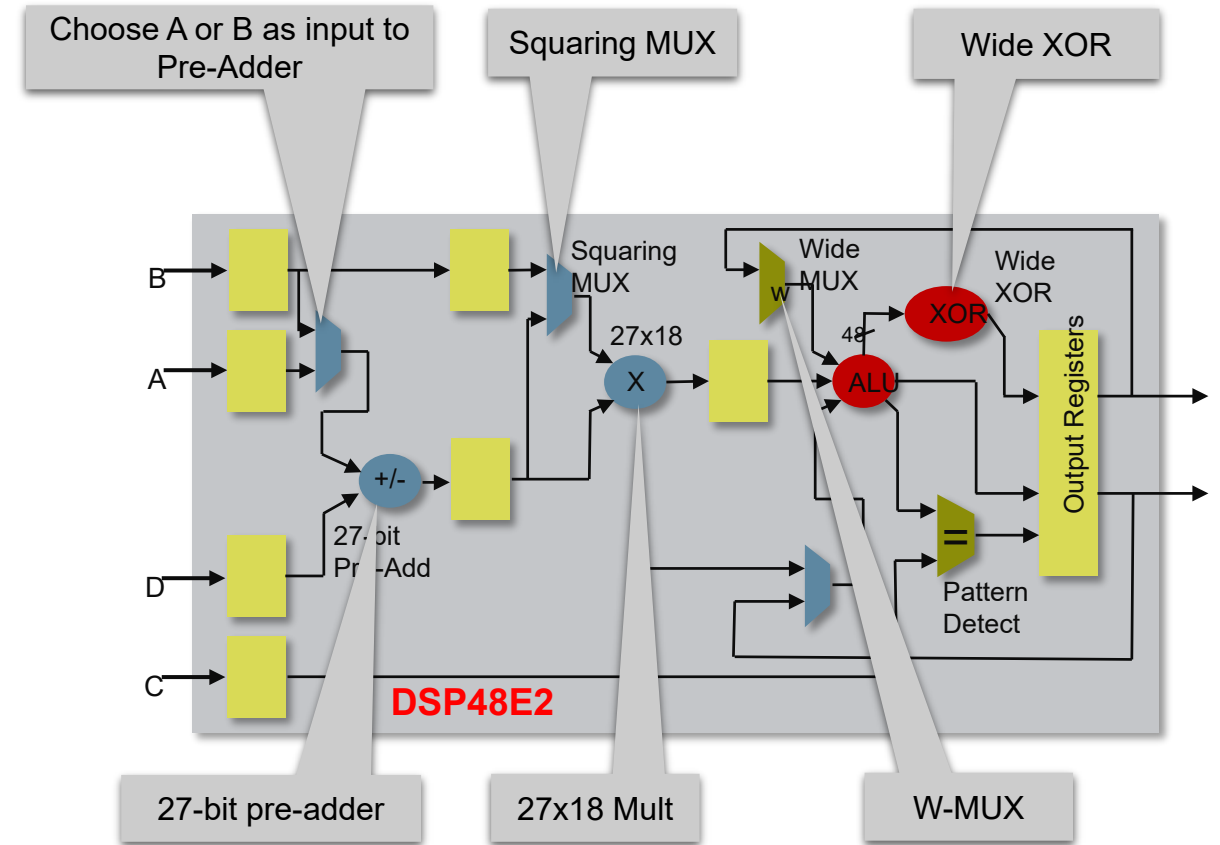


Programmable Routing:

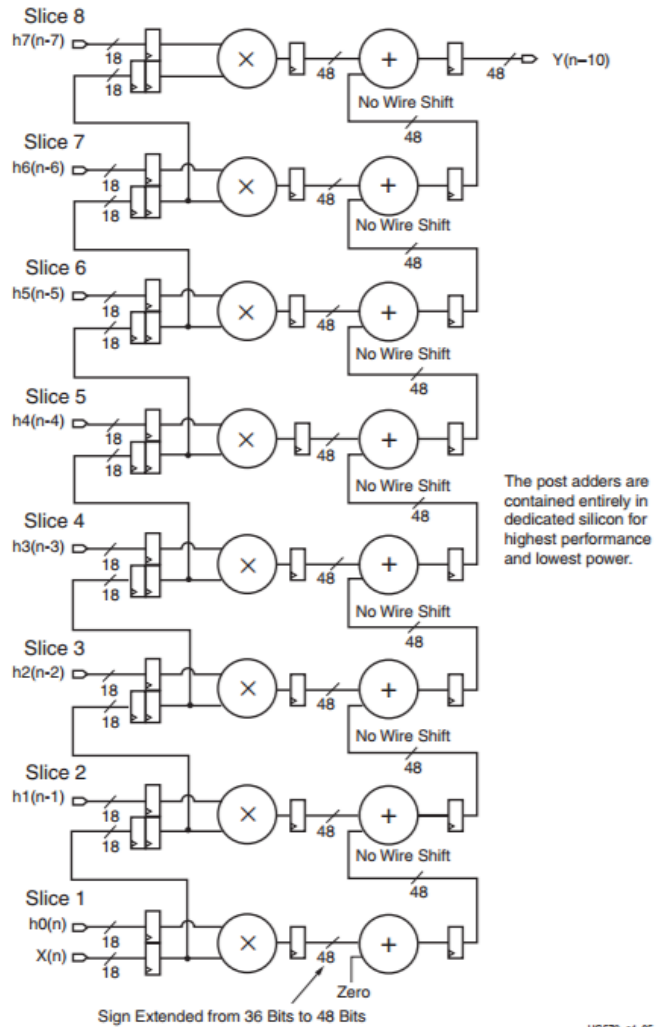
Local routes spanning one,
two, four or five tiles away
Long routes spanning twelve
or sixteen tiles away

DSP48E2 Slice In UltraScale

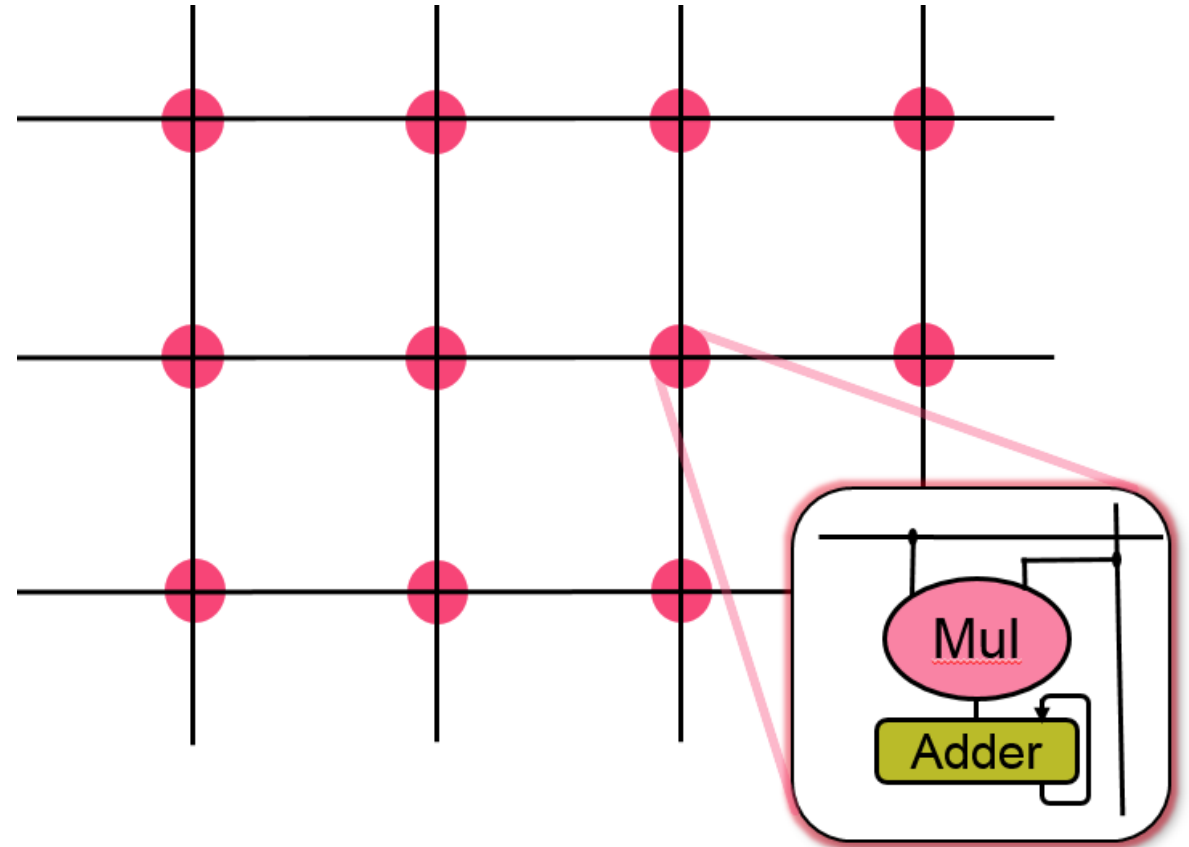
- 27x18-bit multiplier and 27-bit pre-adder
- Pre-adder MUX to select A or B input
- Squaring MUX on output of pre-adder
- Wide XOR (up to 96-bits)
- WMUX product feedback
- Cascading of Input and Outputs



DSP Cascade – Deep learning



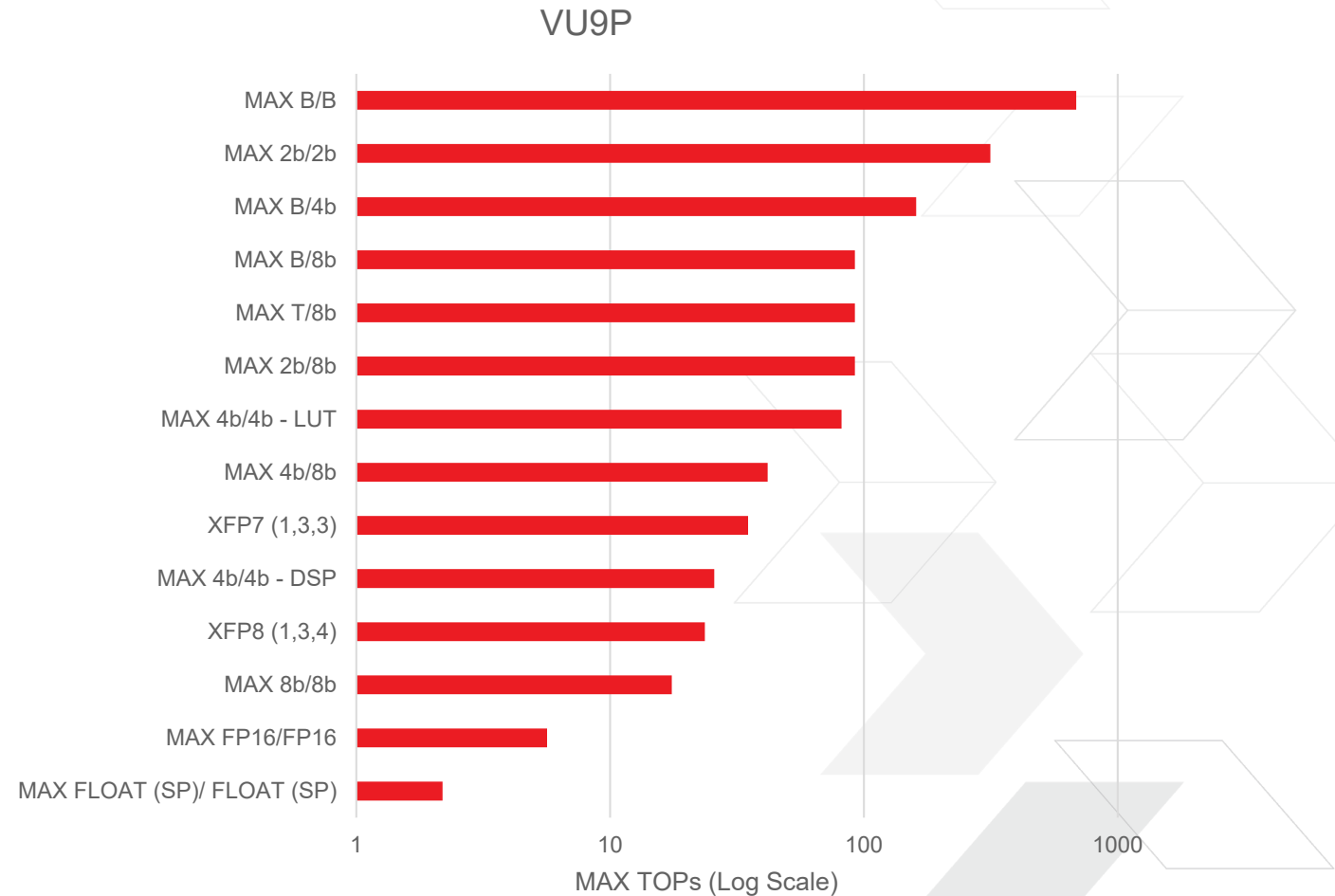
Dot Product Engine – Using Output Cascade



Systolic Array of Accumulators– Using Input Cascade

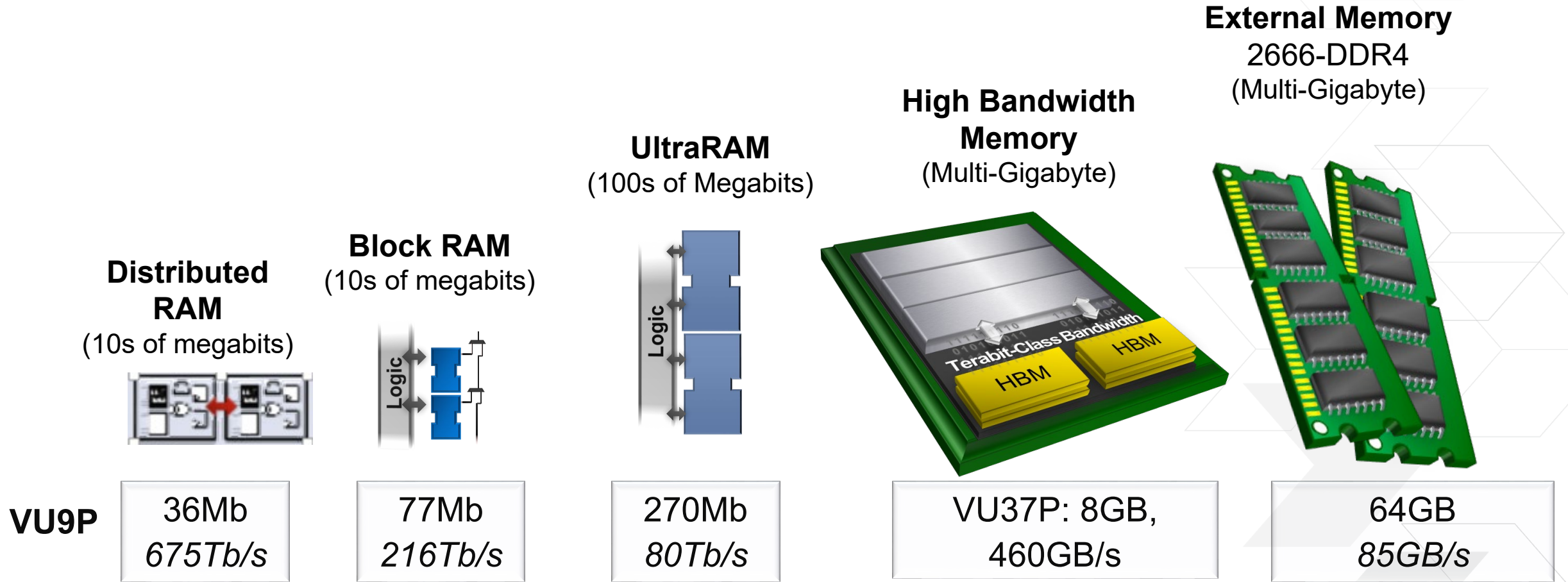
Variable Precision Compute Density – TOPs on VU9P

Weight/Activation	VU9P
MAX FLOAT (SP)/ FLOAT (SP)	2.18538
MAX FP16/FP16	5.64
MAX 8b/8b	17.48304
XFP8 (1,3,4)	23.60306
MAX 4b/4b - DSP	25.71035
XFP7 (1,3,3)	34.96608
MAX 4b/8b	41.72917
MAX 4b/4b - LUT	81.65384
MAX 2b/8b	92.10951
MAX T/8b	92.10951
MAX B/8b	92.10951
MAX B/4b	160.7017
MAX 2b/2b	314.7075
MAX B/B	686.6345



MAX TOPs Estimates at 700 MHz FMAX

Virtex® UltraScale+™ Full Spectrum of Memory

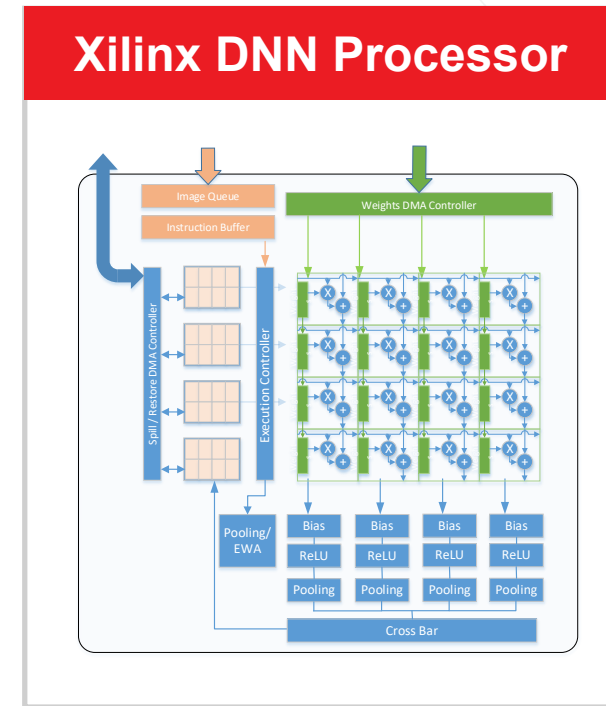
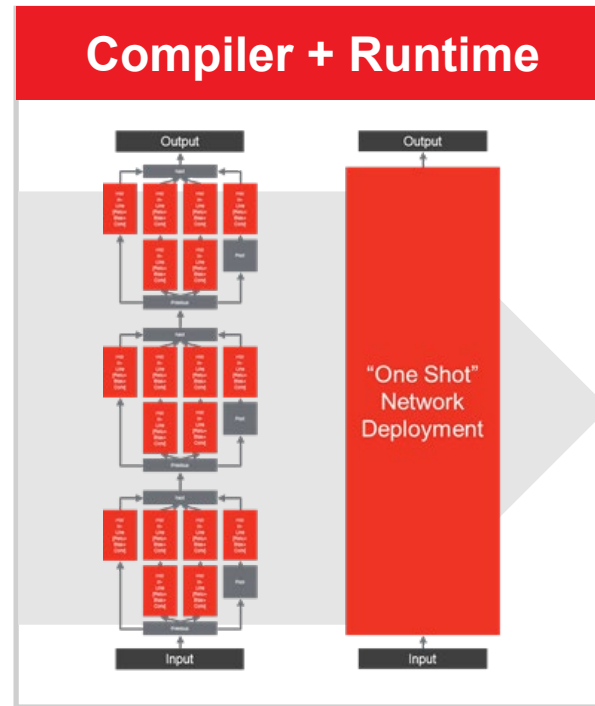
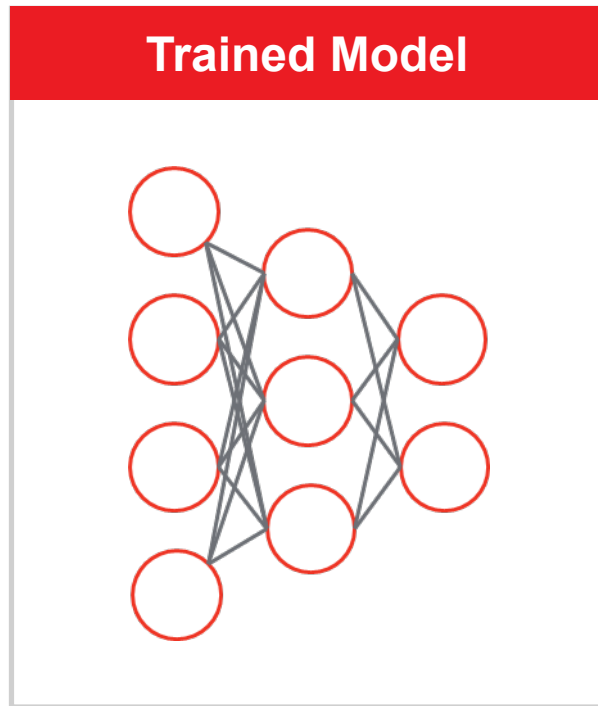


**5 Tiers of Memory -> Build custom memory hierarchy.
500 Mb of On-chip Memory and Tb/s of On-chip Memory Bandwidth**

Xilinx ML-Suite – xDNN Inference Engine

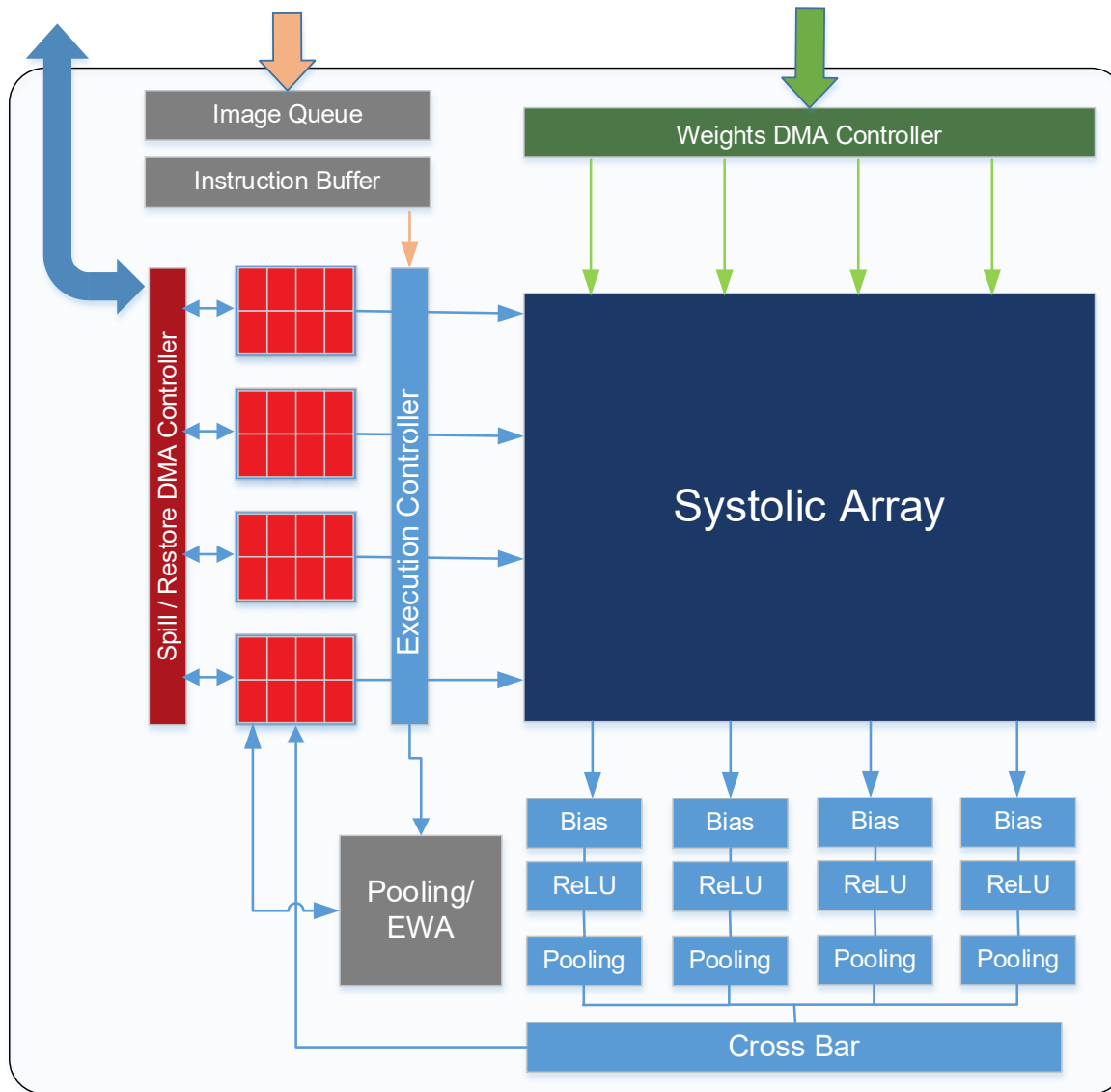


Xilinx Inference Engine – DNN Processor + Compiler



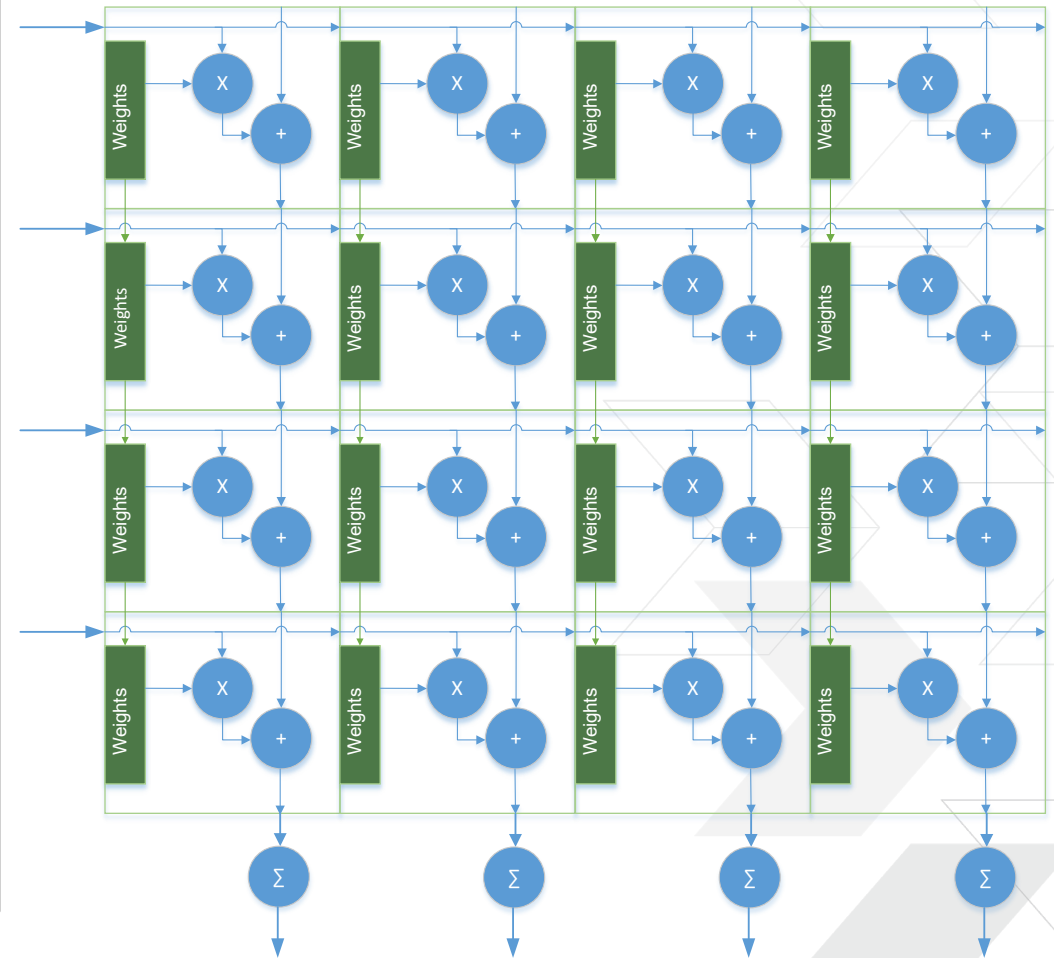
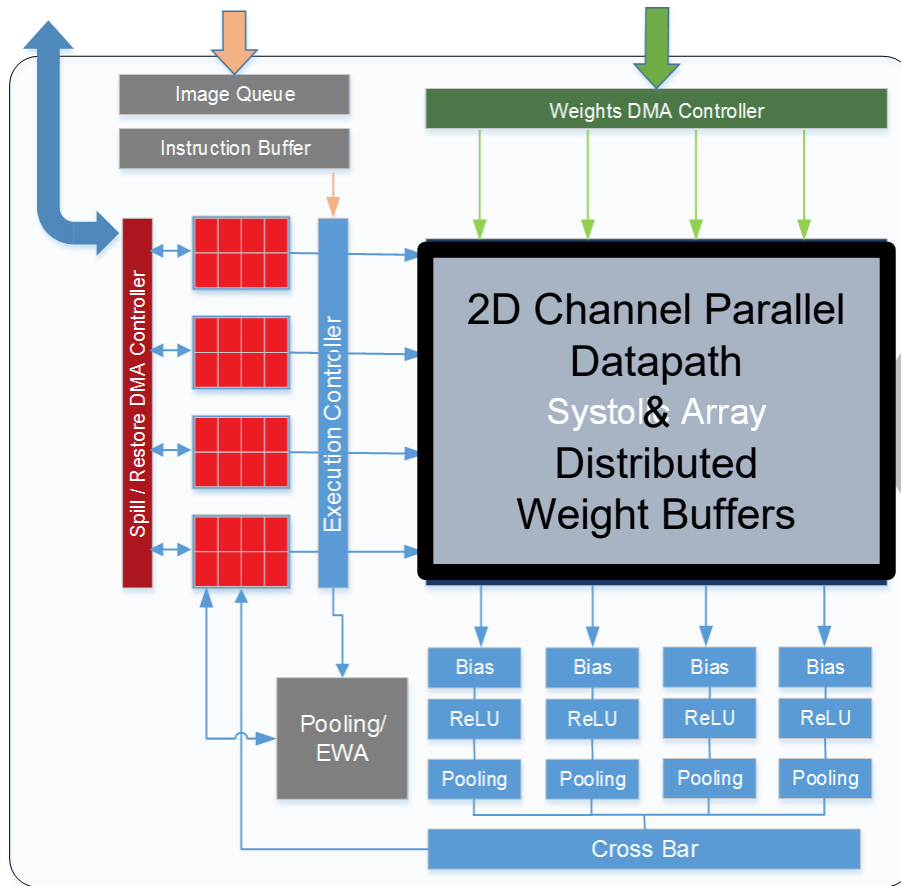
- > Low Latency AND High Throughput at Batch = 1
- > No FPGA Expertise Needed

Xilinx DNN Processor (xDNN)



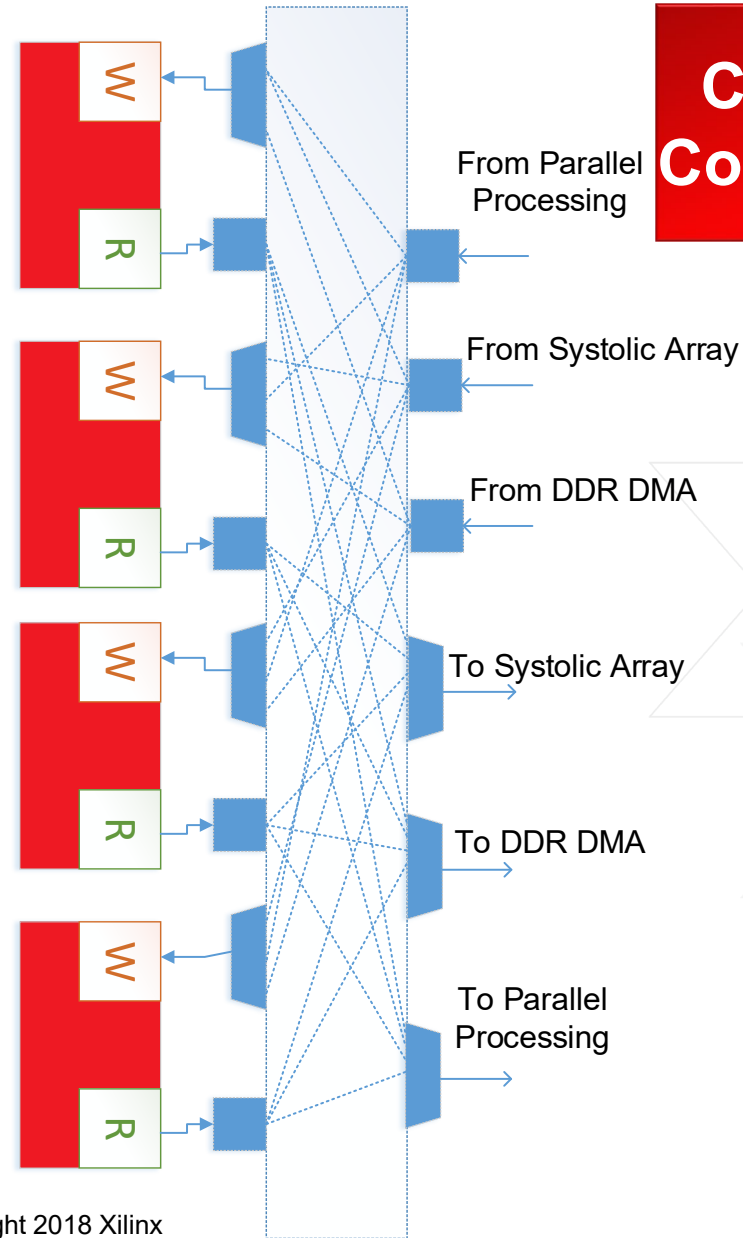
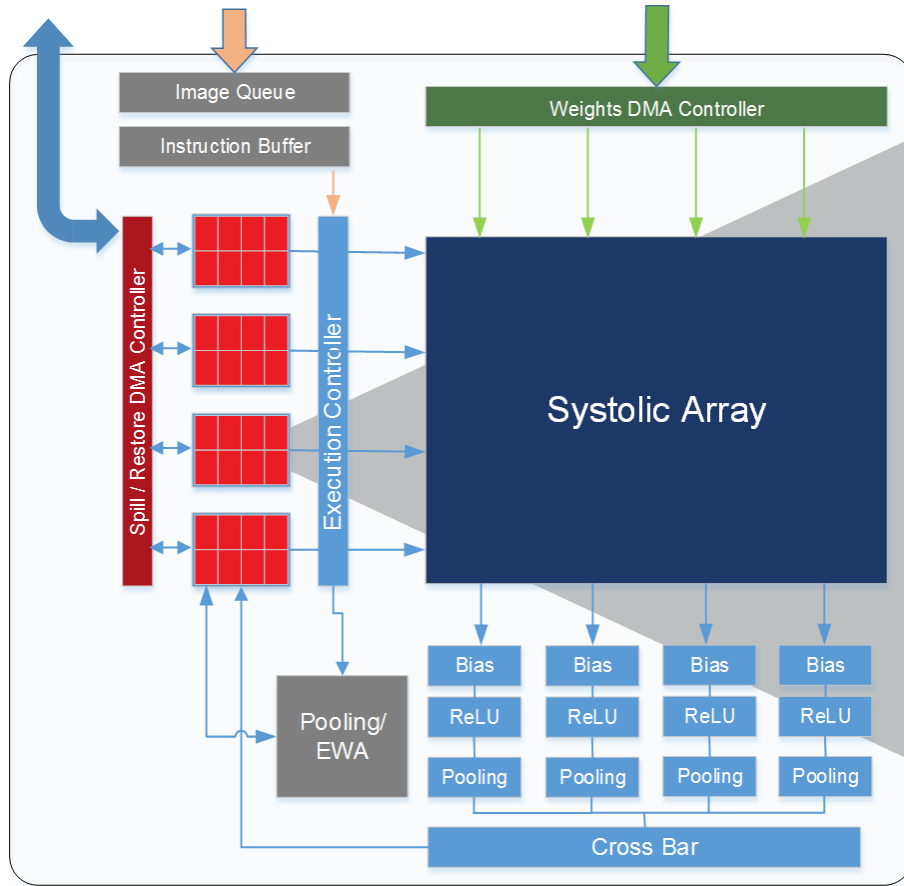
- Low latency CNN Inference engine for Classification, Detection and Segmentation
- Overlay architecture – No need to reprogram FPGA
- Programmable Feature-set
- Tensor Level Instructions
- 700+MHz DSP Freq (VU9P)
- Custom Network Acceleration

xDNN – Channel Parallel Systolic Array



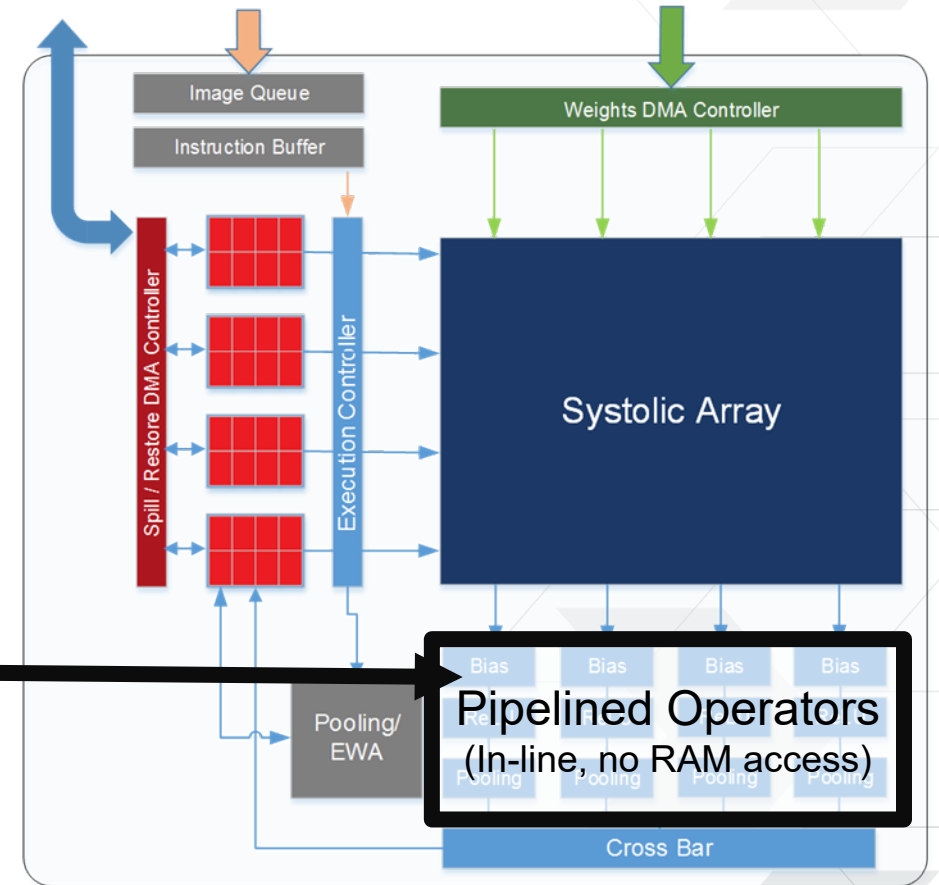
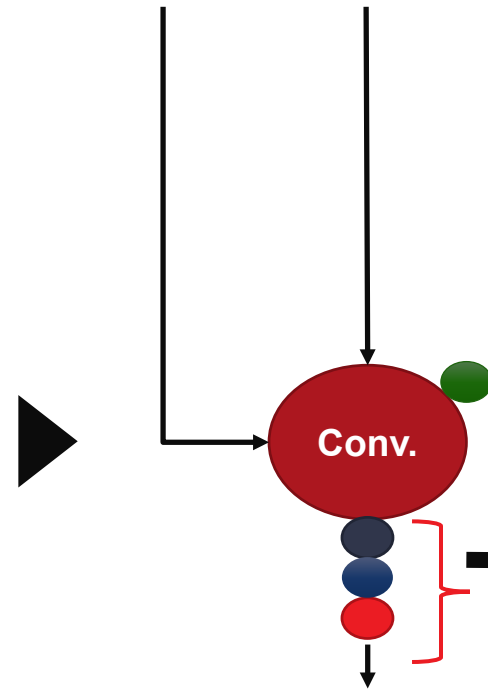
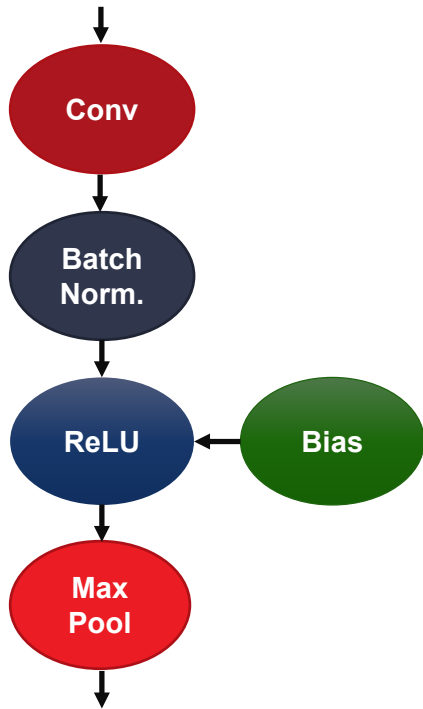
Micro-Architecture Optimized for underlying Ultrascale+ FPGA Fabric

xDNN Processor – Tensor Memory



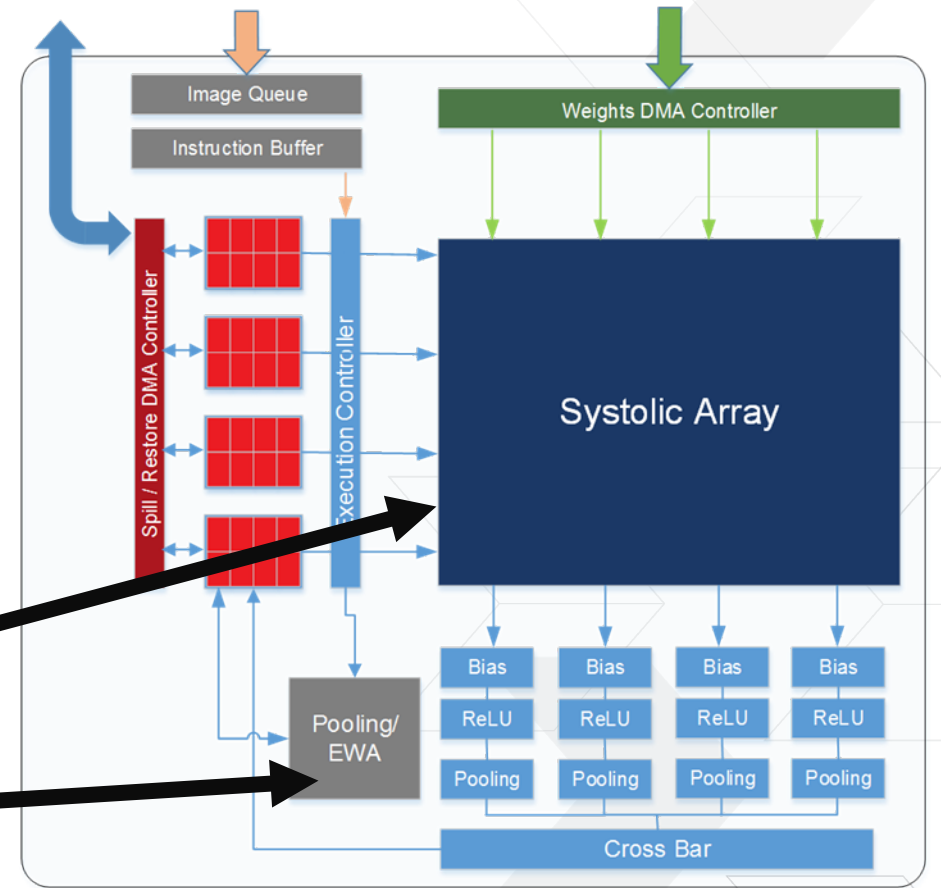
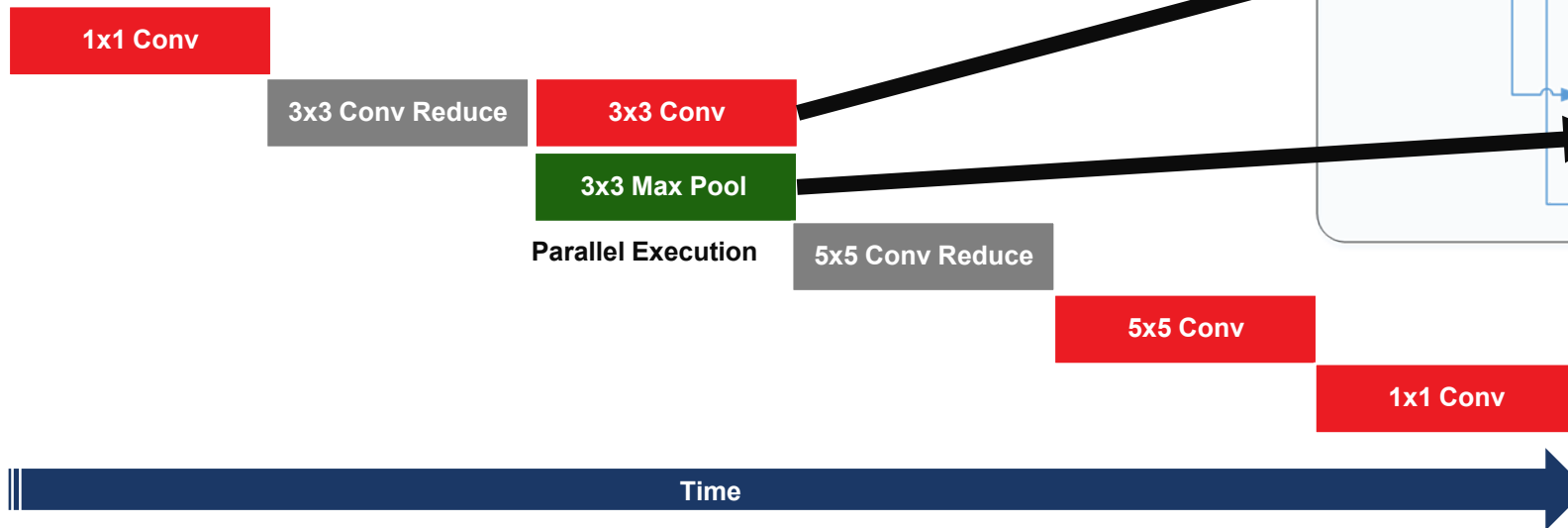
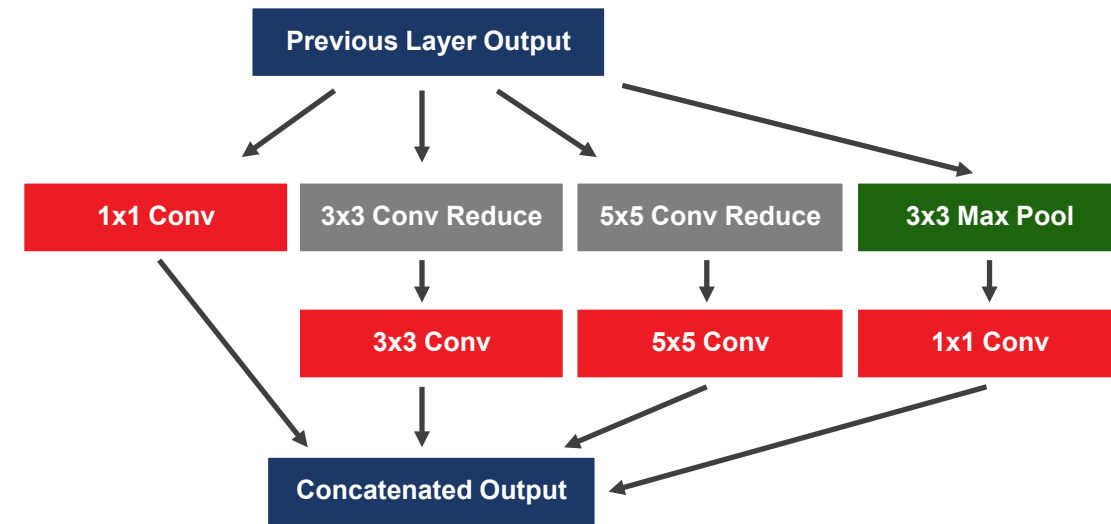
**Channel Parallel
Concurrent Access**

Fusing of Operations



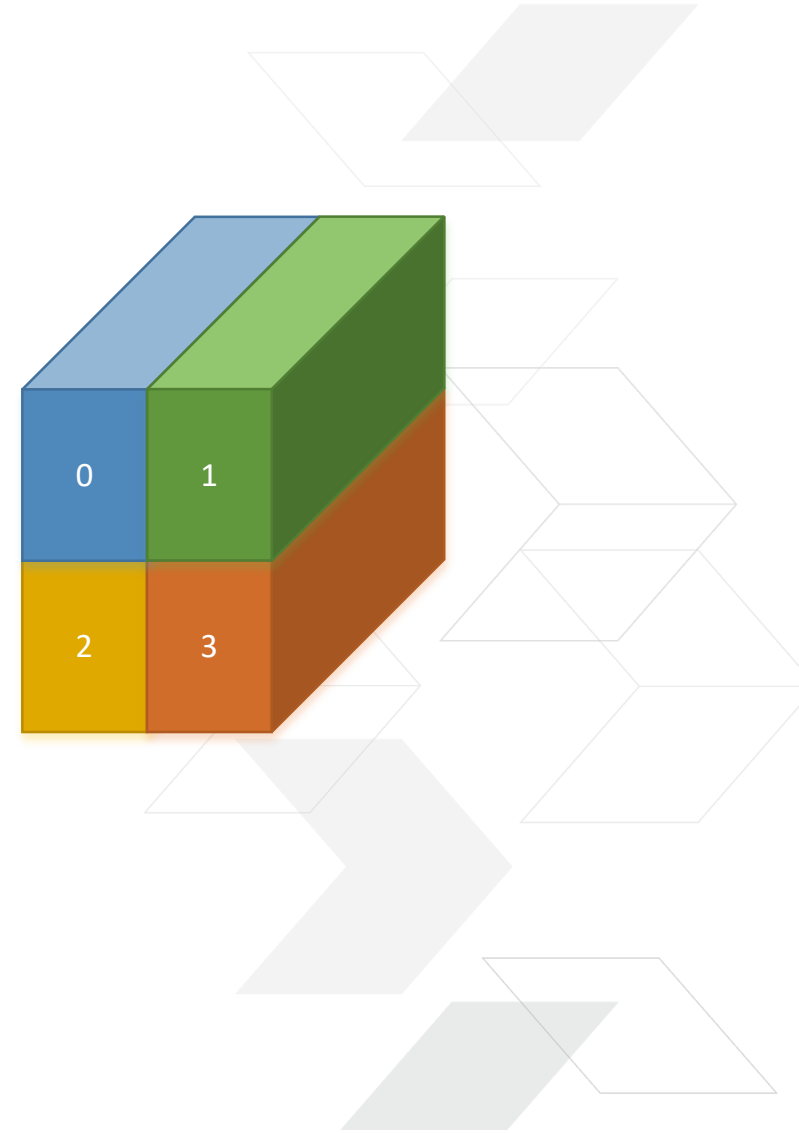
- > Fuse operations as pipe-lined stages
- > Data Flow – No memory access

Instruction Level Parallelism



Automatic Intra Layer Tiling

- **Very important for High Resolution images**
- **Only when feature map size exceeds on-chip memory size**
- **Split each convolution**
 - Break IFM planes to tiles
 - Work on full feature map depth
 - Tiles fetched from external memory
 - Memory re-allocated so tiles overwrite previous tiles



xDNN Feature support

Feature	Details
Convolution	NxM, Stride 1,2,4,8 N,M=1-15
Max Pool	NxM, Stride 1,2,4,8 N,M=1-15
Avg Pool	NxM, Stride 1,2,4,8 N,M=1-15
Dilated Convolution	Factor 1,2,4
De-Convolution	NxM, Stride 1,2,4,8 N,M=1-15
Up-sampling	Factor 2,4,8,16
Activation	ReLU, pReLU
Elementwise Addition	Any – Square, Rectangular
Precision	Int8, Int16
Network	Classification: e.g. ResNet Object Detection: e.g. YOLO v2, Segmentation: e.g. MaskRCNN




Xilinx ML-Suite - xfDNN Software Stack















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

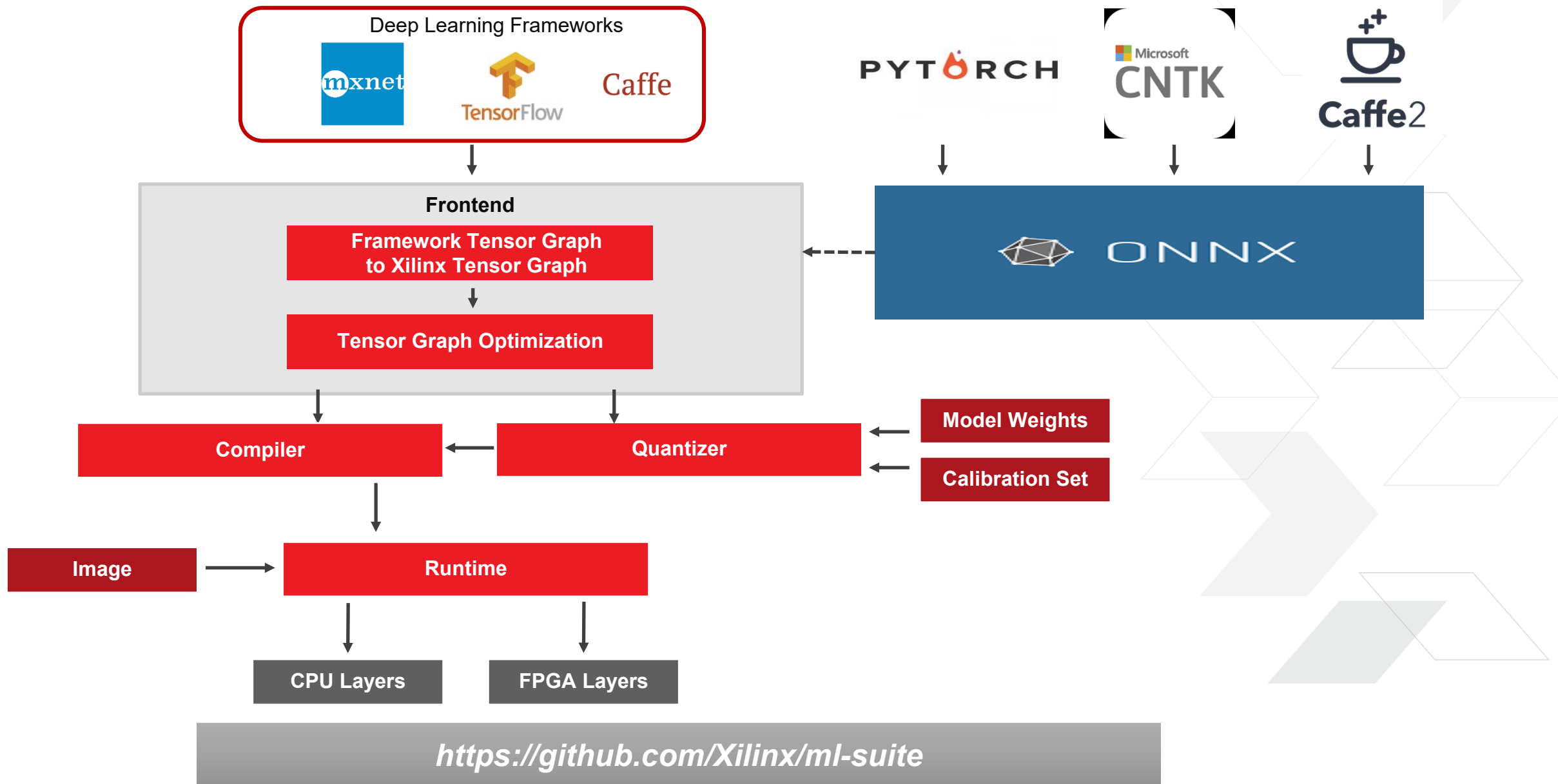
Xilinx ML Suite

Branch: master ▾ [New pull request](#) [Find file](#) [Clone or download ▾](#)

 wilderfield Merge pull request #64 from Xilinx/v1.3 [...](#) Latest commit 5434529 2 days ago

 apps	Adding darknet tools	4 days ago
 docs	Staging 1.3	10 days ago
 examples	Staging 1.3	10 days ago
 ext	Staging 1.3	10 days ago
 notebooks	Updating developer lab	4 days ago
 overlaybins	Add XDNNv3 build for AWS	3 days ago
 tests	Staging 1.3	10 days ago
 xfdnn	Update xfdnn libraries	2 days ago
 .gitignore	Update .gitignore	5 months ago
 LICENSE	Release 1.1 to master	5 months ago
 README.md	Updated docs	8 days ago
 fix_caffe_opencv_symlink.sh	Staging 1.2	2 months ago

xfDNN Flow



xfDNN Compiler

- Caffe - `xfdnn_compiler_caffe.pyc`
- Keras - `xfdnn_compiler_keras.pyc`
- MxNet - `xfdnn_compiler_mxnet.pyc`
- Tensorflow - `xfdnn_compiler_tensorflow.pyc`

Each of these tools contain largely the same arguments. For example, `xfdnn_compiler_caffe` has the following arguments:

```
usage: xfdnn_compiler_caffe.py [-h] [-n NETWORKFILE] [-g GENERATEFILE]
                               [-w WEIGHTS] [-o PNGFILE] [-c CONCATSTRATEGY]
                               [-s STRATEGY] [--schedulefile SCHEDULEFILE]
                               [-i DSP] [-v] [-m MEMORY] [-d DDR] [-p PHASE]
                               [-r RANKDIR]
```

- `-h, --help` - shows this help message and exit
- `-n, --networkfile` - NETWORKFILE Input prototxt for compiler
- `-g, --generatefile` - GENERATEFILE Output of Compiler, xDNN instructions
- `-w, --weights` - WEIGHTS Output of weights for use with deployment APIs
- `-o, --pngfile` - PNGFILE Outputs Optimized Graph as a PNG file, Requires dot executable
- `-c, --concatstrategy` - CONCATSTRATEGY
- `-s, --strategy` - STRATEGY Strategies for compiler (default: all)
- `--schedulefile` - SCHEDULEFILE Show layer by layer memory layout (optional)
- `-i, --dsp` DSP xdnn kernel dimension (28 or 56, default: 28) 28 is also known as "Med" and 56 as "Large"
- `-v, --verbose`
- `-m, --memory MEMORY` - On-chip Memory available in MB (default: 4). Needs to match intended overlaybin.
- `-d, --ddr DDR` - Users DDR memory along with on-chip memory. If the network is too large for the device, set this to 256 .
- `-p, --phase PHASE` - Caffe prototxt phase (TRAIN, TEST, ALL default: TEST)
- `-r, --rankdir RANKDIR`

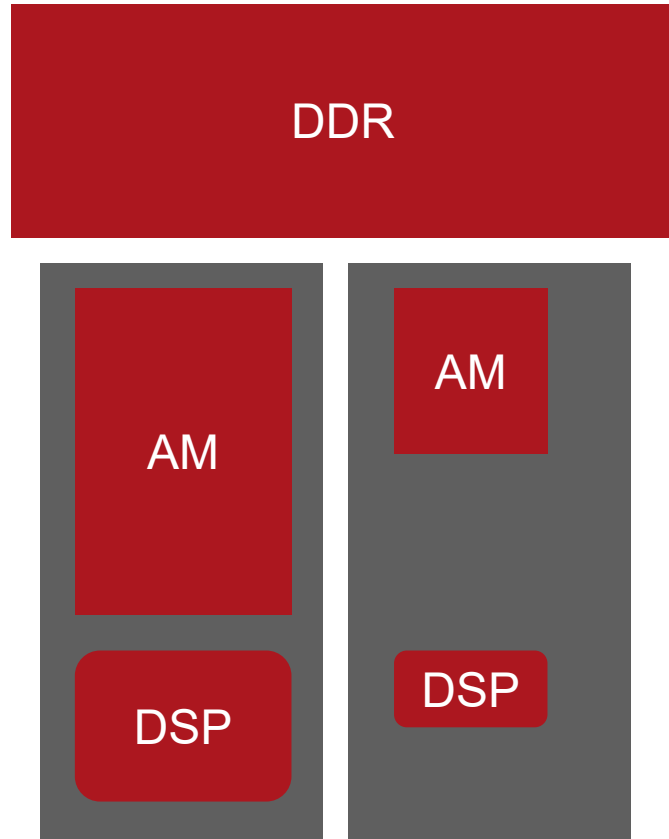
xfDNN Quantizer

```
quantize.pyc [-h] [--deploy_model DEPLOY_MODEL]
              [--output_json OUTPUT_JSON] [--weights WEIGHTS]
              [--calibration_directory CALIBRATION_DIRECTORY]
              [--calibration_size CALIBRATION_SIZE]
              [--calibration_seed CALIBRATION_SEED]
              [--calibration_indices CALIBRATION_INDICES]
              [--bitwidths BITWIDTHS] [--dims DIMS]
              [--transpose TRANSPOSE] [--channel_swap CHANNEL_SWAP]
              [--raw_scale RAW_SCALE] [--mean_value MEAN_VALUE]
              [--input_scale INPUT_SCALE]
```

xfDNN Compiler

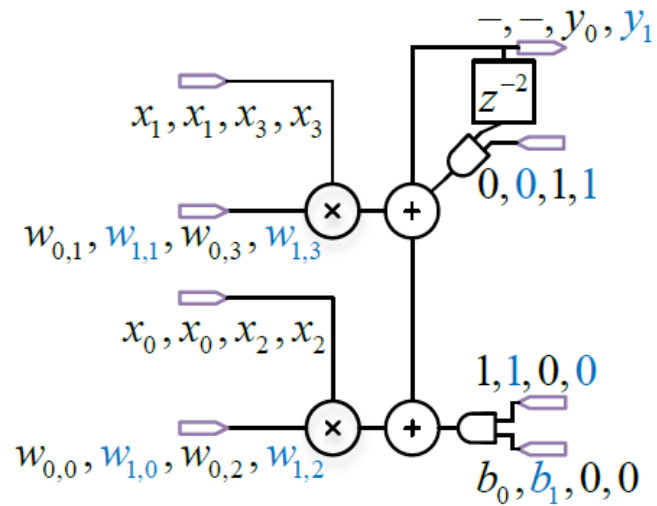
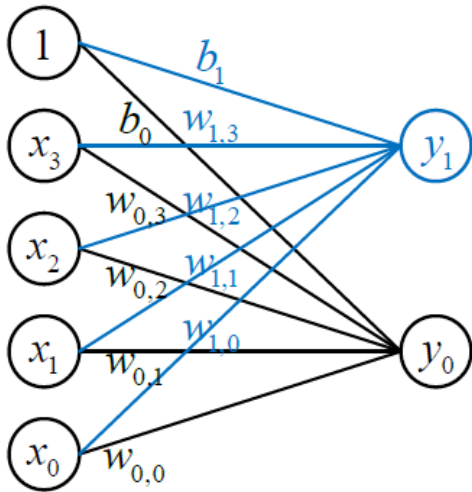


xDNN Logical View

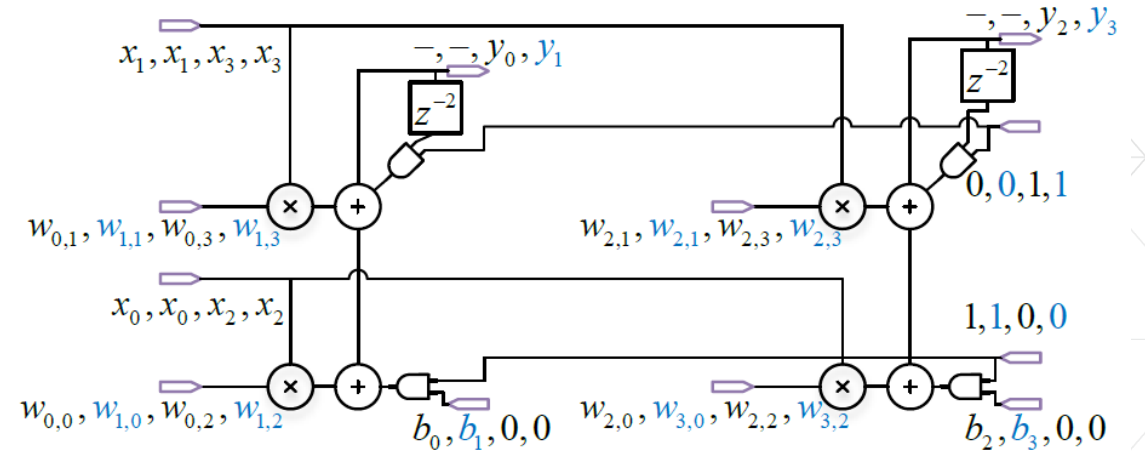


- > **All instructions are high level instructions**
 - >> CONV, POOL, ELEMWISE, ReLU, Shift, Scale
 - >> No need to break down to lower level instructions
- > **DSP array operates on Activation Memory**
 - >> Fast – On Chip SRAM
 - >> Slow – DDR
 - >> No need of explicit data movement during xDNN instructions(V3)
- > **Weights are streamed automatically by xDNN – Formatted by software**
- > **First Activation layer requires software formatting**

DSP - Time Sliced Dot Product Computation



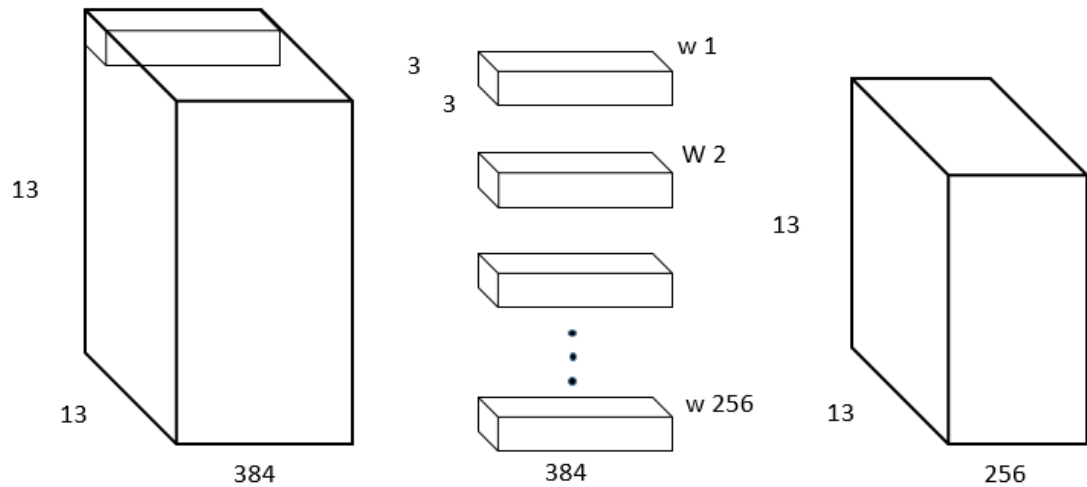
Single Column



Multiple Column

xDNN 96x16 array computes 4 outputs per column

Convolution Operation on xDNN Systolic Array

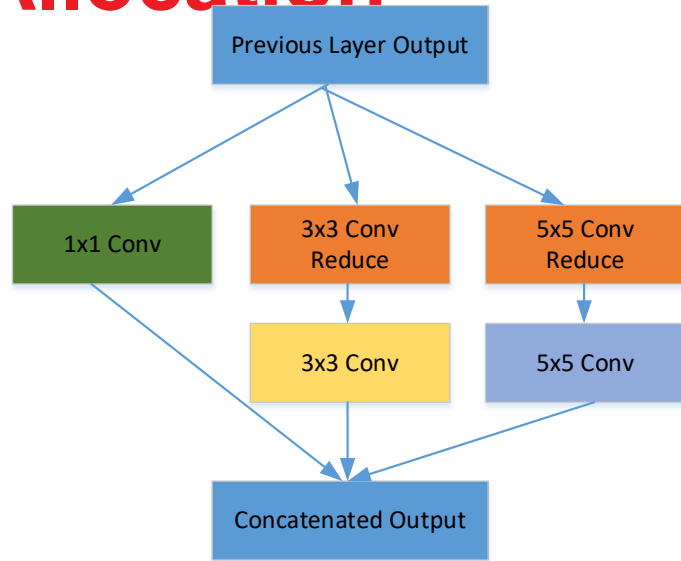


$$y_{v,u,d_2} = \sum_{d_1=0}^{D_1-1} e_{F_{v,u}}^T \left(f^{(l)}(X, v, u, d_1) \odot h^{(l)}(W, v, u, d_1, d_2) \right) e_{H_{v,u}} + b_{d_2}$$

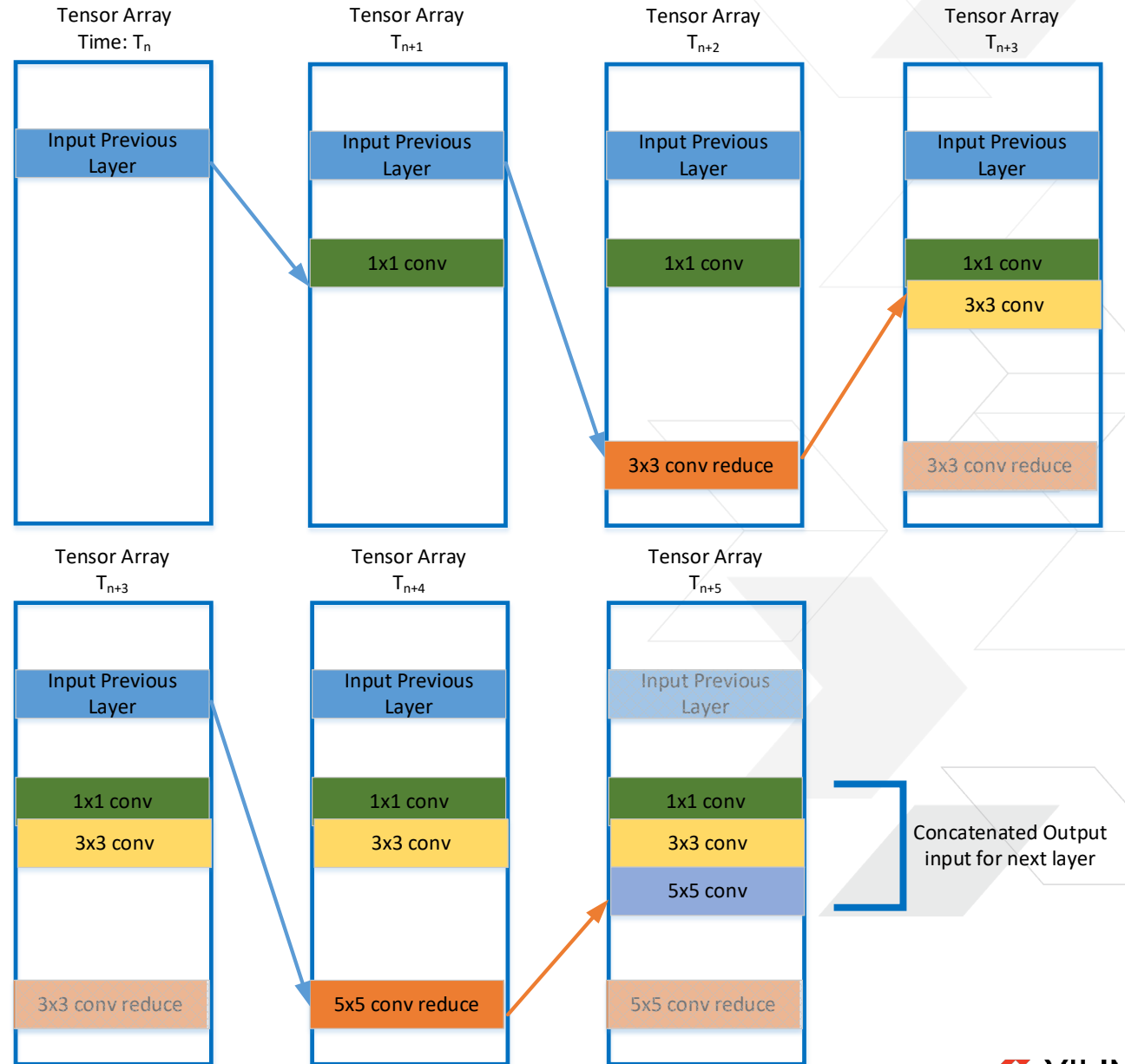
	Bank 0	Bank 1	Bank 2	Bank 3	Bank 4	Bank 5	Bank 6	Bank 7	Bank 8	Bank 9	Bank 10	Bank 11
P[h-1,w-1]												
...												
P[1,w-1]	Ch 192	Ch 200	Ch 208	Ch 216	Ch 224	Ch 232	Ch 240	Ch 248				
P[1,0]												
P[0,w-1]	Ch 199	Ch 207	Ch 215	Ch 223	Ch 231	Ch 239	Ch 247	Ch 255				
...												
P[0,2]												
P[0,1]												
P[0,0]												
...												
P[h-1,w-1]	Ch 96	Ch 104	Ch 112	Ch 120	Ch 128	Ch 136	Ch 144	Ch 152	Ch 160	Ch 168	Ch 176	Ch 184
P[1,1]												
P[1,0]	Ch 103	Ch 111	Ch 119	Ch 127	Ch 135	Ch 143	Ch 151	Ch 159	Ch 167	Ch 175	Ch 183	Ch 191
P[0,w-1]												
...												
P[h-1,w-1]	Ch 0	Ch 8	Ch 16	Ch 24	Ch 32	Ch 40	Ch 48	Ch 56	Ch 64	Ch 72	Ch 80	Ch 88
P[1,w-1]												
P[1,0]	Ch 7	Ch 15	Ch 23	Ch 31	Ch 39	Ch 47	Ch 55	Ch 63	Ch 71	Ch 79	Ch 87	Ch 95
P[0,w-1]												
...												
P[0,2]												
P[0,1]												
P[0,0]												

Real Activation Memory Layout of 256 Channels with kernel size h x w

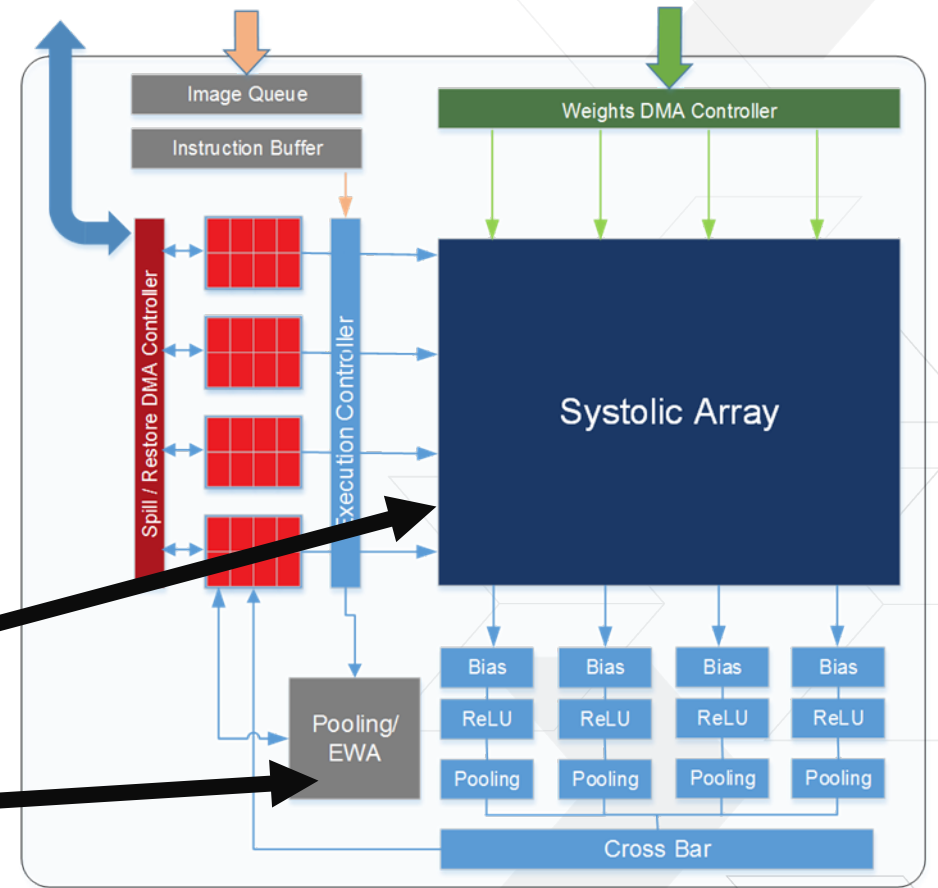
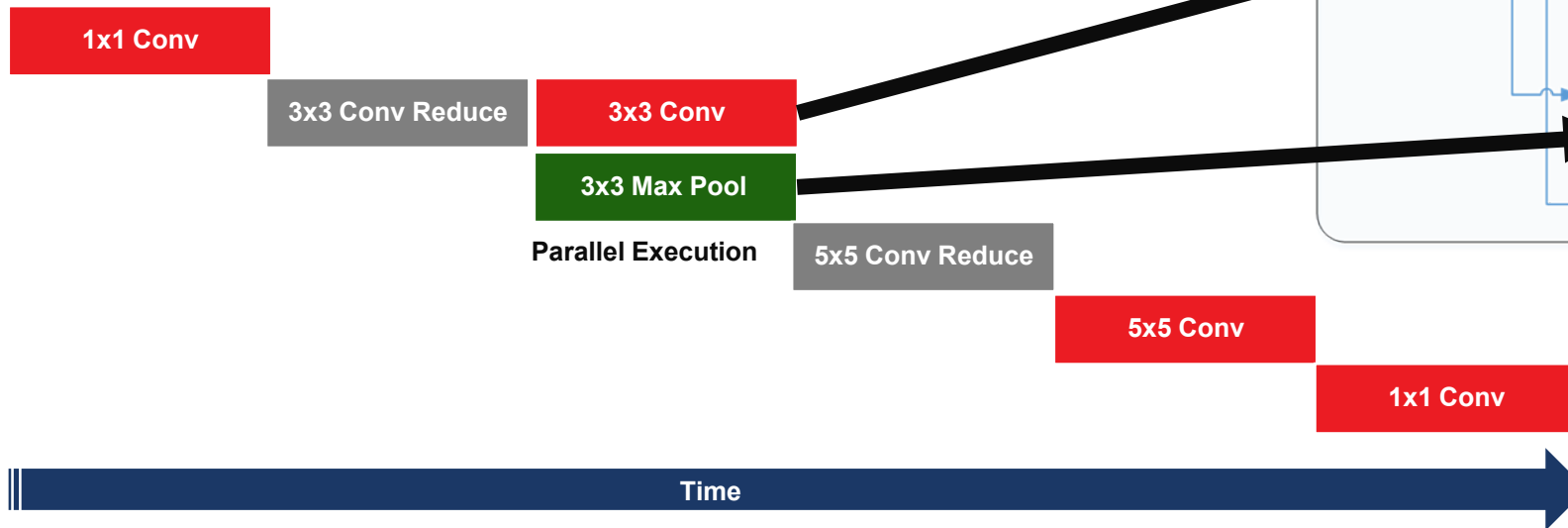
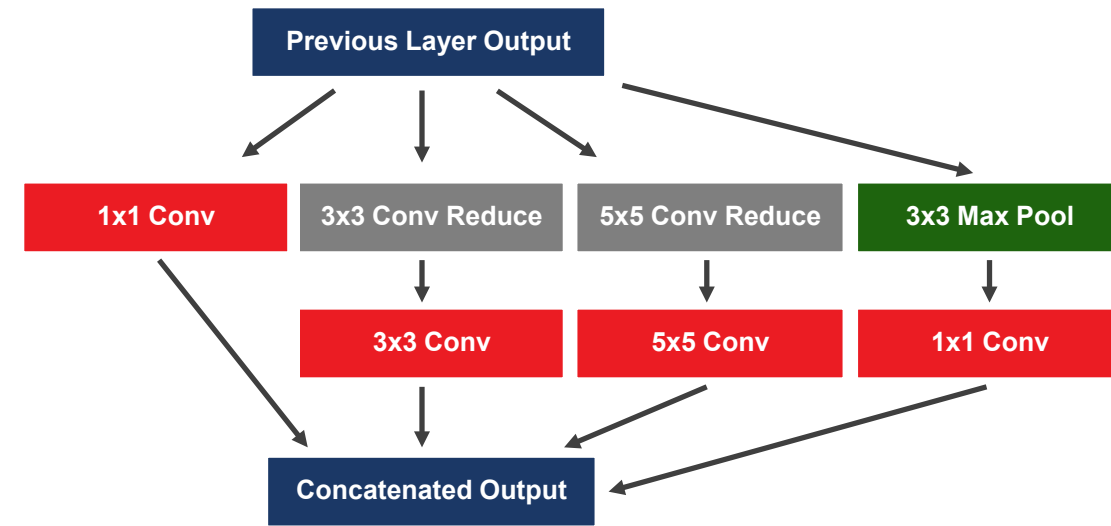
Compiler – Scheduling for Optimized Memory Allocation



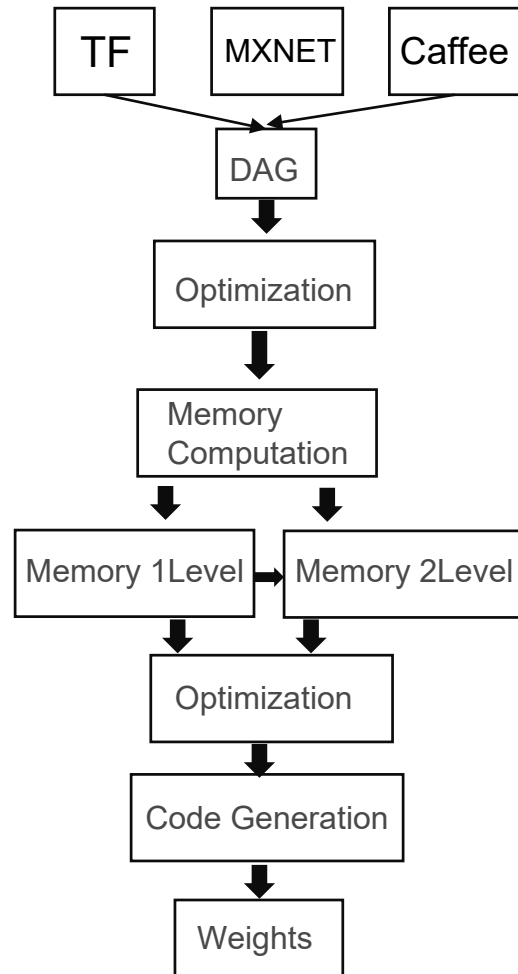
- Address allocation - output of layer n -> input for layer n+1
- Maximize on-chip memory, minimize spill to external memory
- Concatenation through the address generation



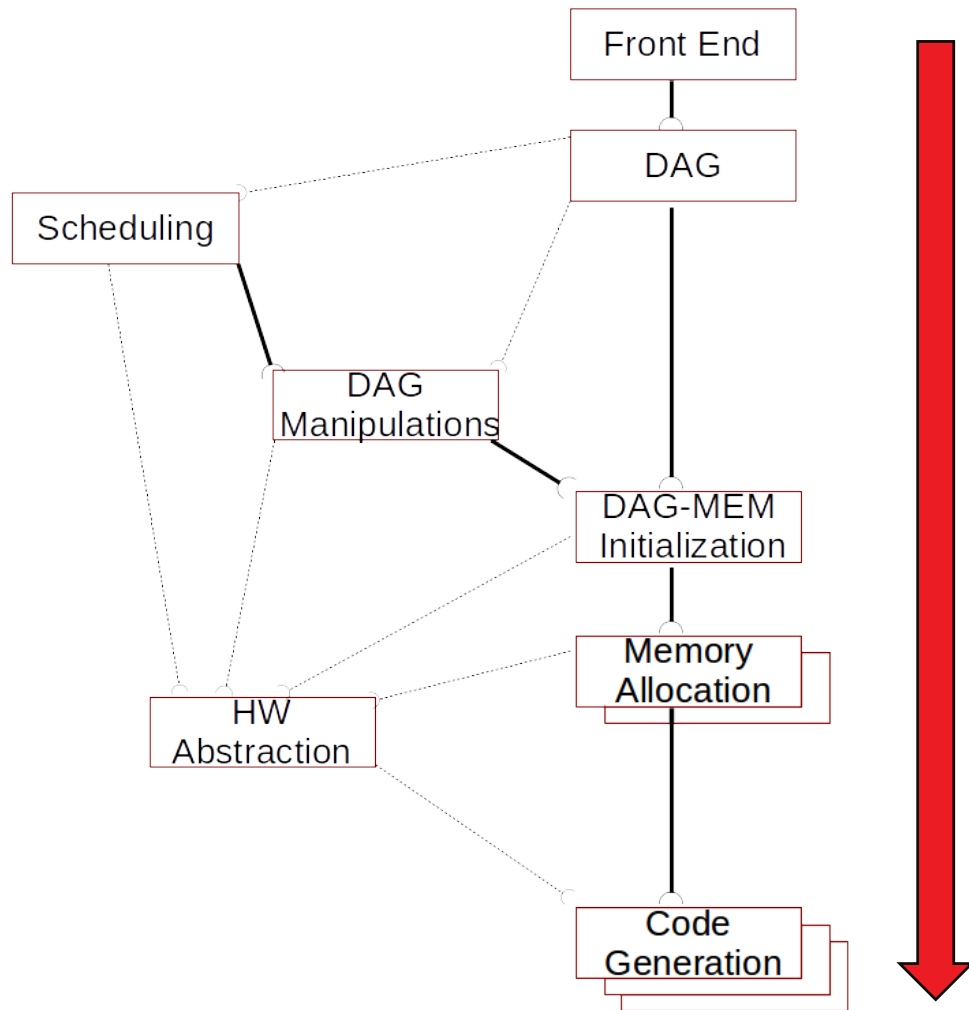
Compiler – Optimize for ILP and Pipelining



xfDNN Compiler Features



- Support for multiple frameworks
 - >> MXNET, TF, Caffe
- Weight extraction and formatting
- DAG optimization and Operator Fusing
 - >> Conv = Conv + BN + Scale
 - >> Relu Merging
 - >> BatchNorm into Scale
 - >> conv_1x1_s2 = conv_1x1 + maxpool
- Slicing support for Conv, Elemwise, Concat
- Support for Scatter and Gather for URAM to DDR data movement and back
- Memory Computation
 - >> Minimum memory and graph initialization
- Memory allocation (Single and 2level)
 - >> URAM – User selectable LRU Cache
 - >> DDR – 256MB
- Code generation
 - >> Data movement + Split Code



The Compiler Organization

Front-end → Code generation (Obviously)

There is interdependency among modules

- Schedules ask for tensors space
- So do DAG manipulations
- Only after Memory allocation
 - we know Tensors space
- At Code generation
 - AM compute or DDR
- Key Optimizations
 - Replication
 - Pooling Around
 - Tiling
 - Parallel Schedule

Replication – Memory Centric and Rule Based

> What is Replication?

- >> The input tensor of a Convolution must satisfy alignment requirements
- >> $R = Channels \% DSP$
 - R is extra space
- >> We can use R to replicate the input tensor in order to exploit parallelism
 - We can double down on R if necessary

> Replication is a manipulation of Tensor layout (space)

> Replication improves the throughput of the convolution unit (time)

- >> If the convolution is not 1x1

> Replication is an example of space-time

Note: Replication is a tensor property. It affects its shapes and its memory space.

	Bank 0	Bank 1	Bank 2	Bank 3	Bank 4	Bank 5	Bank 6	Bank 7	Bank 8	Bank 9	Bank 10	Bank 11
P[0..1]												
P[1..1]												
P[0..0]												
P[0..1]	Ch 192	Ch 200	Ch 208	Ch 216	Ch 224	Ch 232	Ch 240	Ch 248				
P[0..0]	Ch 199	Ch 207	Ch 215	Ch 223	Ch 231	Ch 239	Ch 247	Ch 255				
P[0..1]												
P[1..1]	Ch 96	Ch 104	Ch 112	Ch 120	Ch 128	Ch 136	Ch 144	Ch 152	Ch 160	Ch 168	Ch 176	Ch 184
P[0..0]	Ch 103	Ch 111	Ch 119	Ch 127	Ch 135	Ch 143	Ch 151	Ch 159	Ch 167	Ch 175	Ch 183	Ch 191
P[0..1]												
P[1..1]	Ch 0	Ch 8	Ch 16	Ch 24	Ch 32	Ch 40	Ch 48	Ch 56	Ch 64	Ch 72	Ch 80	Ch 88
P[0..0]	Ch 7	Ch 15	Ch 23	Ch 31	Ch 39	Ch 47	Ch 55	Ch 63	Ch 71	Ch 79	Ch 87	Ch 95

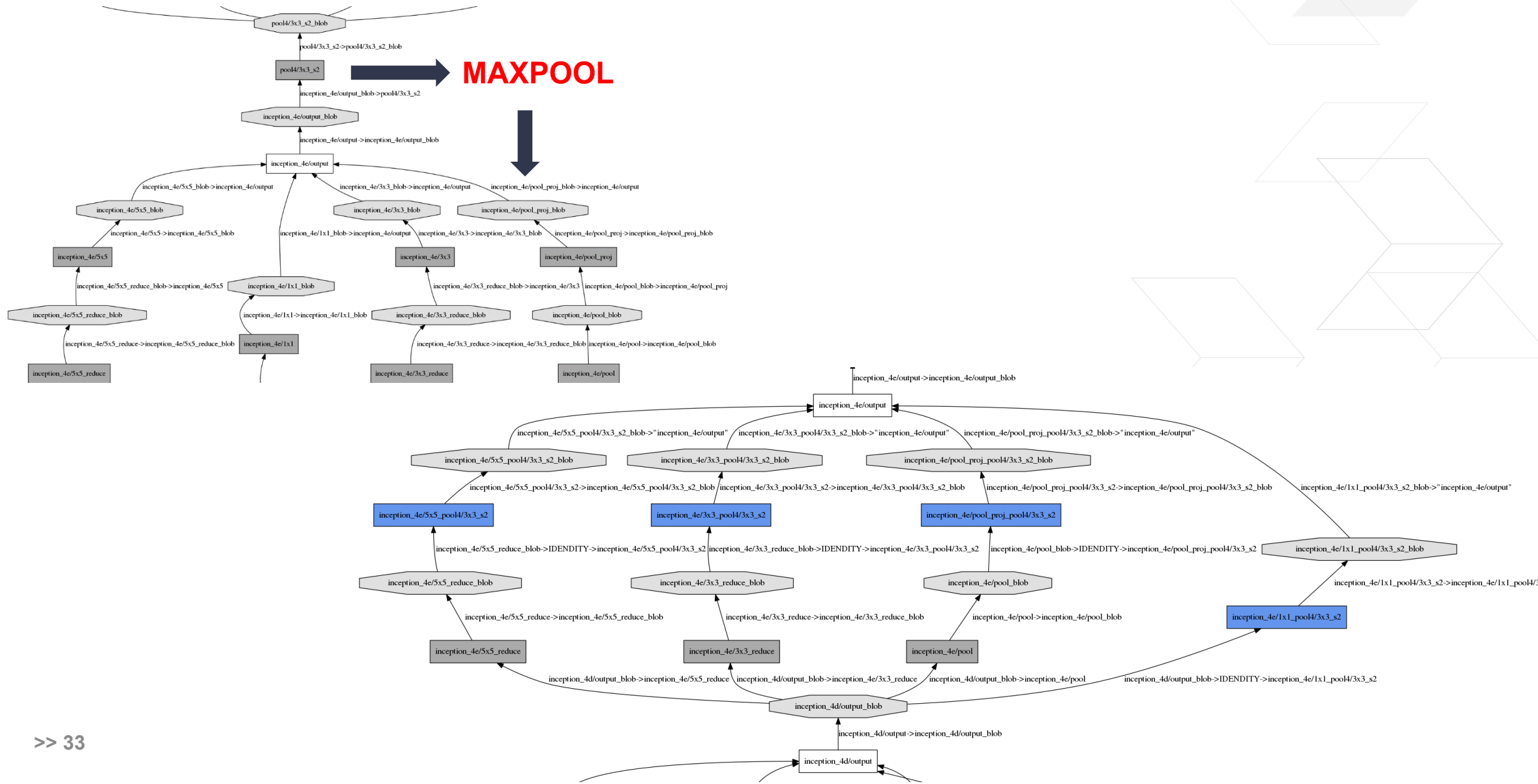
	Bank 0	Bank 1	Bank 2	Bank 3	Bank 4	Bank 5	Bank 6	Bank 7	Bank 8	Bank 9	Bank 10	Bank 11
P[0..1]												
P[1..1]												
P[0..0]												
P[0..1]	Ch 224	Ch 232	Ch 240	Ch 248	Ch 224	Ch 232	Ch 240	Ch 248	Ch 224	Ch 232	Ch 240	Ch 248
P[0..0]	Ch 231	Ch 239	Ch 247	Ch 255	Ch 231	Ch 239	Ch 247	Ch 255	Ch 231	Ch 239	Ch 247	Ch 255
P[0..1]												
P[1..1]	Ch 192	Ch 200	Ch 208	Ch 216	Ch 192	Ch 200	Ch 208	Ch 216	Ch 192	Ch 200	Ch 208	Ch 216
P[0..0]	Ch 199	Ch 207	Ch 215	Ch 223	Ch 199	Ch 207	Ch 215	Ch 223	Ch 199	Ch 207	Ch 215	Ch 223
P[0..1]												
P[1..1]	Ch 96	Ch 104	Ch 112	Ch 120	Ch 128	Ch 136	Ch 144	Ch 152	Ch 160	Ch 168	Ch 176	Ch 184
P[0..0]	Ch 103	Ch 111	Ch 119	Ch 127	Ch 135	Ch 143	Ch 151	Ch 159	Ch 167	Ch 175	Ch 183	Ch 191
P[0..1]												
P[1..1]	Ch 0	Ch 8	Ch 16	Ch 24	Ch 32	Ch 40	Ch 48	Ch 56	Ch 64	Ch 72	Ch 80	Ch 88
P[0..0]	Ch 7	Ch 15	Ch 23	Ch 31	Ch 39	Ch 47	Ch 55	Ch 63	Ch 71	Ch 79	Ch 87	Ch 95

Compiler and Manual Replication

- > **Replication is ON by default**
 - >> If a convolution is not 1x1 it is eligible for Replication
- > **The selection of the details about replication**
 - >> It is hardware dependent
 - >> The compiler uses abstraction
 - And code provided by HW
- > **Replication is a property of the tensors thus:**
 - >> it is applied before memory computation
 - >> And everything looks like any other tensor
- > **We know Replication is important and we allow external inputs**
 - >> Users can overwrite any tensor replication
 - >> By specifying the layer and the input replication



Pooling around ... opportunity and pipelining



Pooling around ... Rule-based graph manipulations

- > **We can move (Max)Pooling across concatenation**
 - >> --poolingaround
- > **We can pipeline Conv + Pooling**
 - >> --pipelineconvmaxpool
- > **This is a two step graph optimization**
 - >> Rule based == always beneficial



Tiling: from DDR to AM and back to DDR

- > **Compiler Memory-Allocation Rules:**
 - >> Either the inputs *and* the outputs are in AM
 - >> Or they are in DDR
- > **Compiler Applies Heuristics**
 - >> Where to allocate: By size, tops, bottoms ..
 - >> How to de-allocate: Least Recently Used policy
- > **Tiling: The data is in DDR**
 - >> V2 the compiler breaks the computation in software tiles
 - AM constraints
 - >> V3 the compiler estimates the largest tile that
 - Minimizes the overlay space/communication
 - Fit in AM
 - >> We do tiling at Code Generation time (last stage)



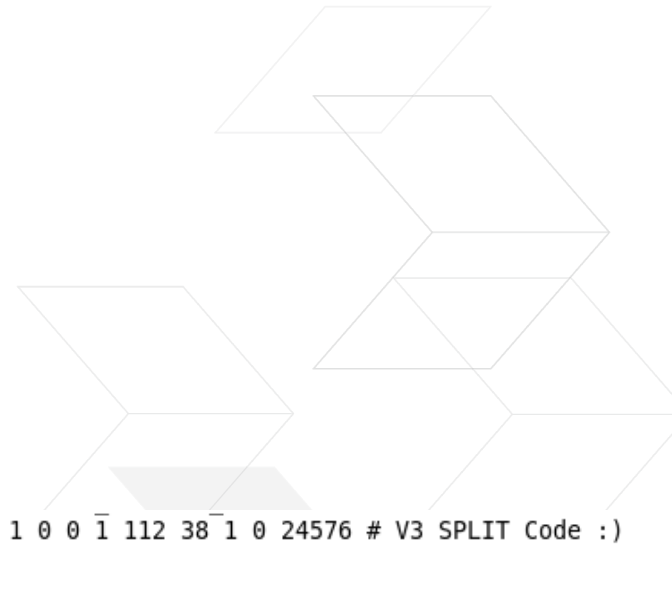
Tiling Examples

V2 can split only the **Height** of the input/output volume

- Channels and width are fixed
 - The minimum tile can be larger than the AM (stop)
 - Multiple instructions

V2 Tiling

```
# spLiT conv1/7x7_s2 2097152 ;
4 XNGather 0x0 0x1c0000 224 224 3 0 1 1 0 113 # # Conv conv1/7x7_s2 # SPLIT Code :)
5 XNConv conv1/7x7_s2#0 7 7 2 2 3 3 1 1 16 26 2 1 1 0x0 224 114 3 0x39000 112 56 64 0 # SPLIT Code :)
6 XNScatter 0x39000 0x0 112 112 64 0 1 1 0 55 # # Conv conv1/7x7_s2 # SPLIT Code :)
7 XNGather 0x0 0x1c0000 224 224 3 0 1 1 109 223 # # Conv conv1/7x7_s2 # SPLIT Code :)
8 XNConv conv1/7x7_s2#1 7 7 2 2 3 0 1 1 16 26 2 1 1 0x0 224 115 3 0x39800 112 56 64 0 # SPLIT Code :)
9 XNScatter 0x39800 0x0 112 112 64 0 1 1 56 111 # # Conv conv1/7x7_s2 # SPLIT Code :)
```



V3 Tiling

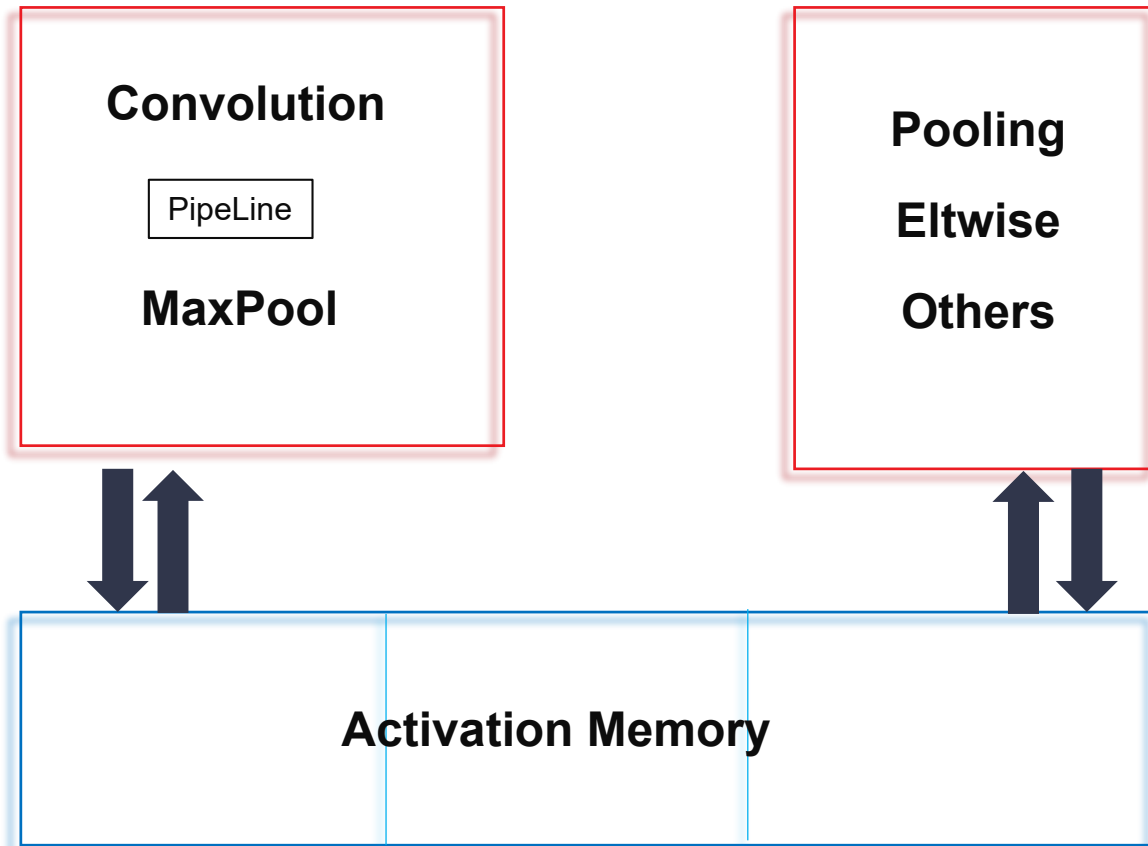
```
3 XNConv conv1/7x7_s2 7 7 2 2 3 3 1 1 16 26 2 1 1 0x0 224 224 3 0x126000 112 112 64 0 0 1 3 32 1 0 0 0 0 0 0 0 0 1 1 1 1 1 0 0 1 0 0 1 112 38 1 0 24576 # V3 SPLIT Code :)
6 XNGather 0xc40 0x126000 112 112 64 0 1 1 0 111 0 # conv1/7x7_s2_blob
7 XNMaxPool pool1/3x3_s2 3 3 2 2 0 0 0xc40 112 112 64 0x0 56 56 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 1 1 1 1 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0
```

V3 can split **Height and Width** of the input/output volume

- Channels are fixed
- The minimum tile is about the size of the channels + overlay
- One instruction

```
"xdnn_instr": "XNConv conv1_1 3 3 1 1 1 1 1 16 26 2 1 1 0x0 480 960 3 0xa8c000 480 960 64 \\  
0 1 0 0 0 2 3 32 0 0 0 0 0 0 1 1 1 1 1 1 0 0 1 0 0 480 49 1 0 24576 # U3 SPLIT Code :)",  
"xdnn_kv": {  
  "HAT": "1", "SRCAM-Buffer_0": "0", "SRCAM-Buffer_1": "24576",  
  "XNOP": "XNConv", "bias": "1", "concat_full_sect_num": "0", "concat_i_am_your_father": "0", "concat_repl_sect_num": "0",  
  "concat_repl_unit_num": "0", "concat_repl_unit_width": "0", "concat_starting_ch": "0", "dilation_h": "1", "dilation_w": "1",  
  "dst_full_sect_num": "0", "dst_repl_sect_num": "2", "dst_repl_unit_num": "3", "dst_repl_unit_width": "32", "en_halfmode": "0",  
  "en_inlinemaxpool": "0", "en_pingpong_weight": "1", "id": "2", "inaddr": "0x0", "inchan": "3", "insize_h": "960", "insize_w": "480",  
  "kernel_h": "3", "kernel_w": "3", "name": "conv1_1", "outaddr": "0xa8c000", "outchan": "64", "outside_h": "960", "outside_w": "480",  
  "padding_h": "1", "padding_w": "1", "parallel_read": "0", "postshift": "2", "prerelu": "0", "preshift": "16", "relu": "1",  
  "scale": "26", "slice": "0", "src_full_sect_num": "1", "src_repl_sect_num": "0", "src_repl_unit_num": "0", "src_repl_unit_width": "0",  
  "strides_h": "1", "strides_w": "1", "tile_height": "49", "tile_width": "480", "wait_conv": "1", "wait_download": "1", "wait_ew": "1",  
  "wait_pool": "1", "wait_upload": "1", "wait_upsmp1": "1"  
},
```

Convolution || Pooling ... We have the technology



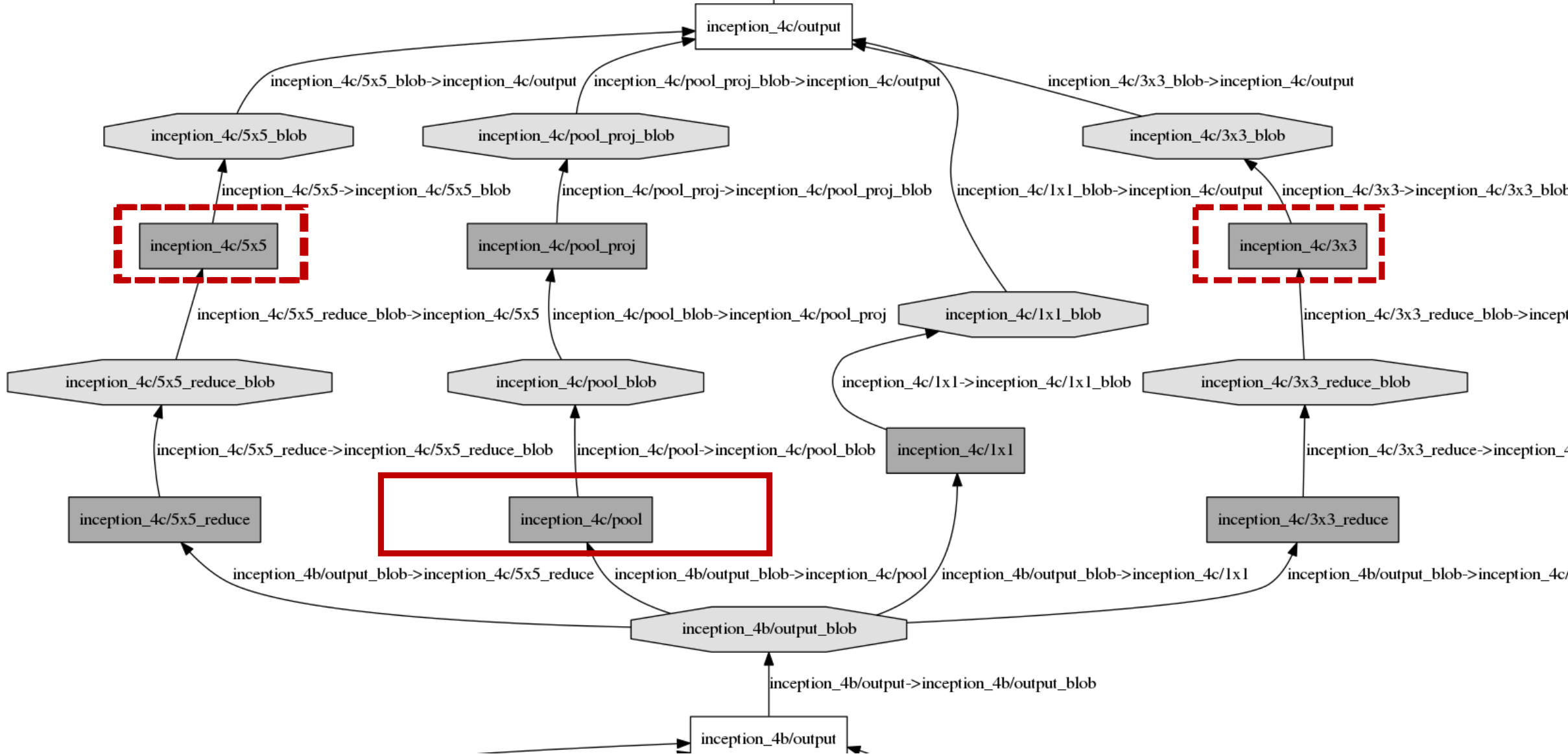
Obviously,

- Two separated kernels
- Parallel access to AM

Translate to,

- Two *close-by* instructions
- No shared Tensors
- And **far apart** Memory Locations

Convolution || Pooling ... Opportunity



Convolution || Pooling ... Codes (V2 simplified)

Regular Schedule: **NOT** parallel and **NOT** optimal

```
41 XNConv inception_4b/pool_proj 1 1 1 1 0 0 1 1 16 26 2 1 1 0x0 14 14 512 0x504000 14 14 64 0
# 42 XNConcat inception_4b/output 0x440000 0x520000 917504
43 XNConv inception_4c/1x1 1 1 1 1 0 0 1 1 16 26 2 1 1 0x440000 14 14 512 0x520000 14 14 128 0
44 XNConv inception_4c/3x3_reduce 1 1 1 1 0 0 1 1 16 26 2 1 1 0x440000 14 14 512 0x0 14 14 128 0
45 XNConv inception_4c/3x3 3 3 1 1 1 1 1 1 16 26 2 1 1 0x0 14 14 128 0x558000 14 14 256 0
46 XNConv inception_4c/5x5_reduce 1 1 1 1 0 0 1 1 16 26 2 1 1 0x440000 14 14 512 0x0 14 14 24 0
47 XNConv inception_4c/5x5 5 5 1 1 2 2 1 1 16 26 2 1 1 0x0 14 14 24 0x5c8000 14 14 64 0
48 XNMaxPool inception_4c/pool 3 3 1 1 1 1 1 0x440000 14 14 512 0x0 14 14 0
49 XNConv inception_4c/pool_proj 1 1 1 1 0 0 1 1 16 26 2 1 1 0x0 14 14 512 0x5e4000 14 14 64 0
# 50 XNConcat inception_4c/output 0x520000 0x600000 917504
51 XNConv inception_4d/1x1 1 1 1 1 0 0 1 1 16 26 2 1 1 0x520000 14 14 512 0x439000 14 14 112 0
```

Optimal Schedule: 😊

Optimality (later)

Data parallelism (short)

Pool

- IN address **High**
- OUT address **Low**

Conv

- IN address **Low**
- OUT address **High**

```
34 XNConv inception_4b/1x1 1 1 1 1 0 0 1 1 16 26 2 1 1 0x36e000 14 14 512 0x520000 14 14 160 0
# 35 XNConcat inception_4b/output 0x520000 0x600000 917504
36 XNConv inception_4c/5x5_reduce 1 1 1 1 0 0 1 1 16 26 2 1 1 0x520000 14 14 512 0x0 14 14 24 0
37 XNConv inception_4c/5x5 5 5 1 1 2 2 1 1 16 26 2 1 1 0x0 14 14 24 0x4e8000 14 14 64 0
38 XNConv inception_4c/3x3_reduce 1 1 1 1 0 0 1 1 16 26 2 1 1 0x520000 14 14 512 0x0 14 14 128 0
39 XNMaxPool inception_4c/pool 3 3 1 1 1 1 1 0x520000 14 14 512 0x38000 14 14 0
39 XNConv inception_4c/3x3 3 3 1 1 1 1 1 1 16 26 2 1 1 0x0 14 14 128 0x478000 14 14 256 0
40 XNConv inception_4c/pool_proj 1 1 1 1 0 0 1 1 16 26 2 1 1 0x38000 14 14 512 0x504000 14 14 64 0
41 XNConv inception_4c/1x1 1 1 1 1 0 0 1 1 16 26 2 1 1 0x520000 14 14 512 0x440000 14 14 128 0
# 42 XNConcat inception_4c/output 0x440000 0x520000 917504
43 XNConv inception_4d/5x5_reduce 1 1 1 1 0 0 1 1 16 26 2 1 1 0x440000 14 14 512 0x520000 14 14 32 0
```


Convolution || Pooling ... Inceptionality ?



An **Inception** is a sub graph

- N(x) specific topological order
- An inception is specified by N
 - Graph property

Inception is self contained

- We explore 6,000 schedules
- Only 600 are valid
 - Up to a max 10,000

Each optimal schedule glued together into a single schedule

```
Inception inception_4b/output - inception_4c/output
depth tfs dfs name [ inputs name']
# 0 18 18 42 inception_4b/output ops 0
#   INPUT [u'inception_4b/5x5', u'inception_4b/1x1', u'inception_4b/pool_proj', u'inception_4b/3x3']
#   OUTPUT [u'inception_4c/1x1', u'inception_4c/pool', u'inception_4c/5x5_reduce', u'inception_4c/3x3_reduce']
# 1 19 19 44 inception_4c/pool ops 1806336
#   INPUT [u'inception_4b/output']
#   OUTPUT [u'inception_4c/pool_proj']
# 2 19 19 46 inception_4c/5x5_reduce ops 9408
#   INPUT [u'inception_4b/output']
#   OUTPUT [u'inception_4c/5x5']
# 3 19 19 43 inception_4c/1x1 ops 50176
#   INPUT [u'inception_4b/output']
#   OUTPUT [u'inception_4c/output']
# 4 19 19 48 inception_4c/3x3_reduce ops 50176
#   INPUT [u'inception_4b/output']
#   OUTPUT [u'inception_4c/3x3']
# 5 20 20 49 inception_4c/3x3 ops 903168
#   INPUT [u'inception_4c/3x3_reduce']
#   OUTPUT [u'inception_4c/output']
# 6 20 20 47 inception_4c/5x5 ops 627200
#   INPUT [u'inception_4c/5x5_reduce']
#   OUTPUT [u'inception_4c/output']
# 7 20 20 45 inception_4c/pool_proj ops 25088
#   INPUT [u'inception_4c/pool']
#   OUTPUT [u'inception_4c/output']
sorted path by memory use
(114688.0, [u'inception_4b/output', u'inception_4c/5x5_reduce', u'inception_4c/5x5', u'inception_4c/output'])
(229376.0, [u'inception_4b/output', u'inception_4c/1x1', u'inception_4c/output'])
(458752.0, [u'inception_4b/output', u'inception_4c/3x3_reduce', u'inception_4c/3x3', u'inception_4c/output'])
(917504.0, [u'inception_4b/output', u'inception_4c/pool', u'inception_4c/pool_proj', u'inception_4c/output'])
5040
tick search 0 of 5040
Number of valid schedules 630
```


Convolution || Pooling ... optimality by Pareto search

The best schedule, with min average space, and min average ops 8 1745408.0 3.75 2568384

```
0 inception_4b/output
1 inception_4c/5x5_reduce
2 inception_4c/5x5
3 inception_4c/3x3_reduce
4 inception_4c/pool
5 inception_4c/3x3
6 inception_4c/pool_proj
7 inception_4c/1x1
8 inception_4c/output
```

Simplified schedule

Committed Ins # 8

```
Concat ▲ I:size 917504 POOLING:None,CONVOLUTION:None,ANY_:inception_4b/output
live:[u'inception_4b/output'] ops:0
▲ I:size 960512 POOLING:None,CONVOLUTION:inception_4c/5x5_reduce,ANY_:None
live:[u'inception_4c/5x5_reduce', u'inception_4b/output'] ops:9408
▲ I:size 1075200 POOLING:None,CONVOLUTION:inception_4c/5x5,ANY_:None
live:[u'inception_4c/5x5_reduce', u'inception_4b/output', u'inception_4c/5x5'] ops:627200
▲ I:size 1261568 POOLING:None,CONVOLUTION:inception_4c/3x3_reduce,ANY_:None
live:[u'inception_4c/3x3_reduce', u'inception_4b/output', u'inception_4c/5x5'] ops:50176
▲ I:size 2637824 POOLING:inception_4c/pool,CONVOLUTION:inception_4c/3x3,ANY_:None
BINGO ▲ live:[u'inception_4c/3x3_reduce', u'inception_4b/output', u'inception_4c/5x5', u'inception_4c/3x3', u'inception_4c/pool'] ops:1806336
I:size 2523136 POOLING:None,CONVOLUTION:inception_4c/pool_proj,ANY_:None
▲ live:[u'inception_4c/3x3', u'inception_4c/pool', u'inception_4c/pool_proj', u'inception_4b/output', u'inception_4c/5x5'] ops:25088
I:size 1835008 POOLING:None,CONVOLUTION:inception_4c/1x1,ANY_:None
▲ live:[u'inception_4c/3x3', u'inception_4c/pool_proj', u'inception_4c/1x1', u'inception_4b/output', u'inception_4c/5x5'] ops:50176
I:size 2752512 POOLING:None,CONVOLUTION:None,ANY_:inception_4c/output
Concat ▲ live:[u'inception_4c/3x3', u'inception_4c/pool_proj', u'inception_4c/1x1', u'inception_4b/output', u'inception_4c/5x5', u'inception_4c/output'] ops:0
The best live 9 1640675.5555555555 3.6666666666666666
0 inception_4b/output
```

Dimensions to minimize:

1. Number of parallel steps
2. Number of live tensors

This is a full search

Convolution || Pooling ... Code generation (V3)

```
PARALLEL >  
inception_4c/pool inception_4c/3x3 Wait(Wait_Download=1, Wait_Upload=1, Wait_Conv=1, Wait_Pool=0, Wait_EW=1, Wait_Upsmpl=1, ParalleRead=0)  
inception_4c/3x3 IN 9361920-9437184 OUT 313600-413952  
inception_4c/pool IN 37632-263424 OUT 9136128-9361920  
-----
```

- > **AS the last step**
 - >> We reach the convolution instruction
 - >> We know that there is a possible parallel Pool
 - >> We check the data parallelism
- > **Parallel**
 - >> Parallel: do not wait for the completion of the Pooling

NOTES: V3 allows explicit parallelism

- > **scheduling, inst-parallelism, and data-parallelism are architecture independent.**

Improvements Due to Compiler Optimizations

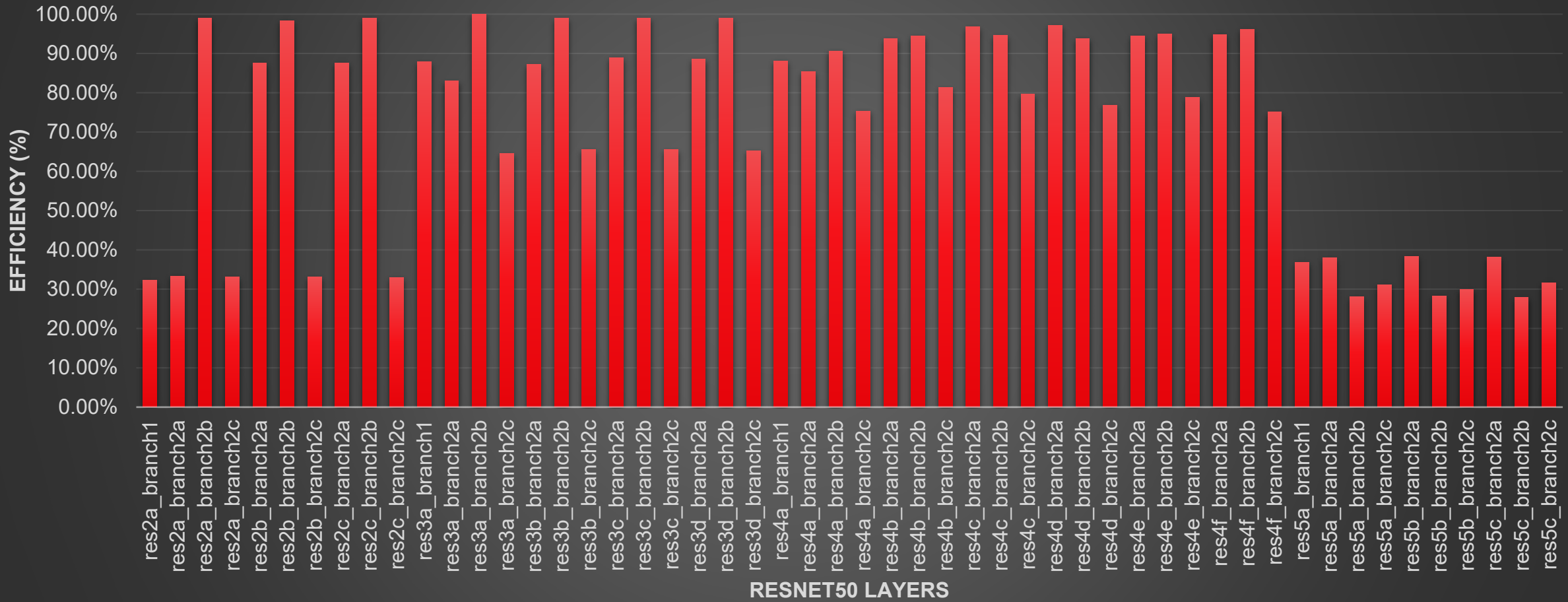
No opt	100.00%
Replication	57.50%
Replication + Conv/Pool pipeline	54.20%
Replication + Conv/Pool pipeline + Parallel	53.97%
Replication + Conv/Pool pipeline + Pooling around	53.18%
Parallel	51.92%
All auto-opts	46.71%
Conv/Pool pipeline + Parallel + dedicated 1 st block	38.68%
All auto-opts + dedicated 1 st block	36.85%
Conv/Pool pipeline + Parallel + “manual” replication	36.69%
All auto-opts + “manual” replication	34.84%

GoogleNetV1

No opt	100.00%
Replication	73.49%
Replication + Conv/Pool pipeline	72.00%
Replication + Conv/Pool pipeline + Parallel	68.88%
Replication + Conv/Pool pipeline + Parallel + “manual” replication	61.01%
All auto-opts + “manual” replication	55.13%

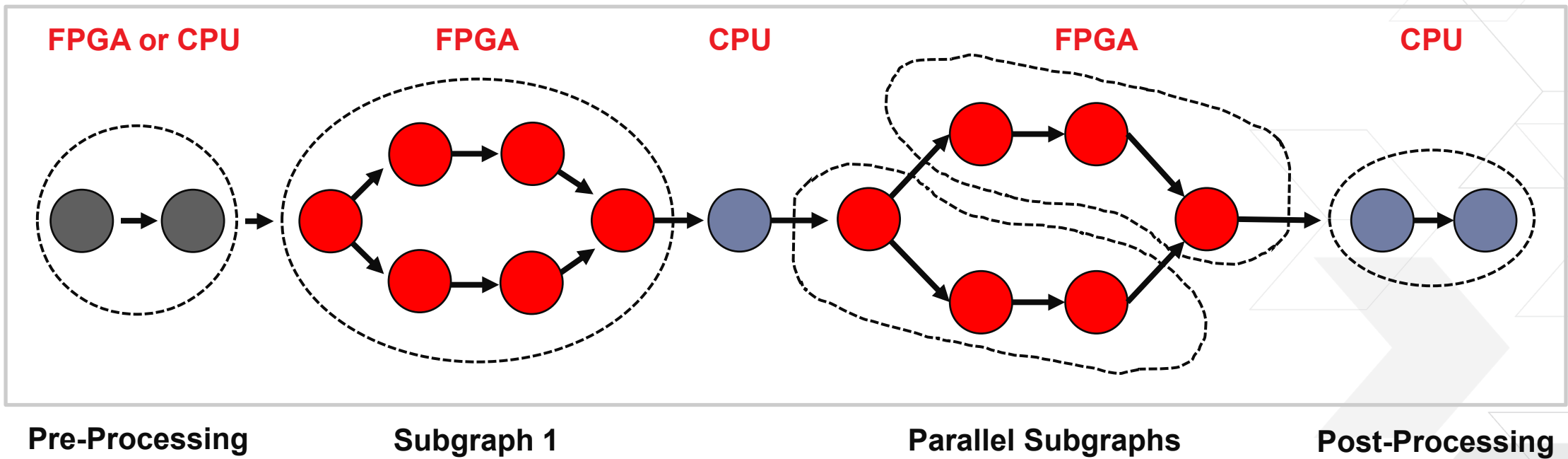
Resnet50

ResNet50 Layer-wise Efficiency



Graph Partitioning

- > Automatically Partition and Code generation of Tensorgraph
- > Data Flow Execution using multiple xDNN/FPGAs



Pre-Processing

Subgraph 1

Parallel Subgraphs

Post-Processing

1 Core -> Multi-Core -> Multi-Chip

xfDNN Quantizer



Quantization

- Quantization:
 - Represent parameter and activation with lower bit width
- Advantages:
 - Smaller
 - Faster
 - Power efficiency
 - Hardware Friendly
- Challenge:
 - Accuracy

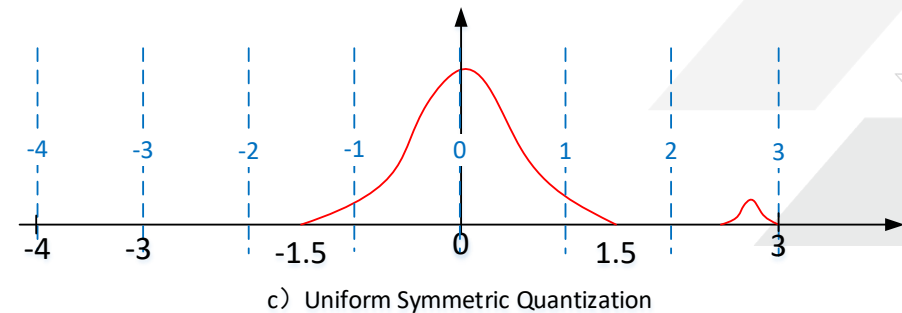
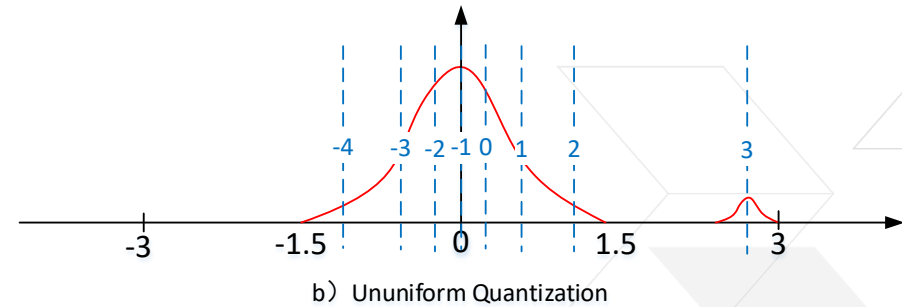
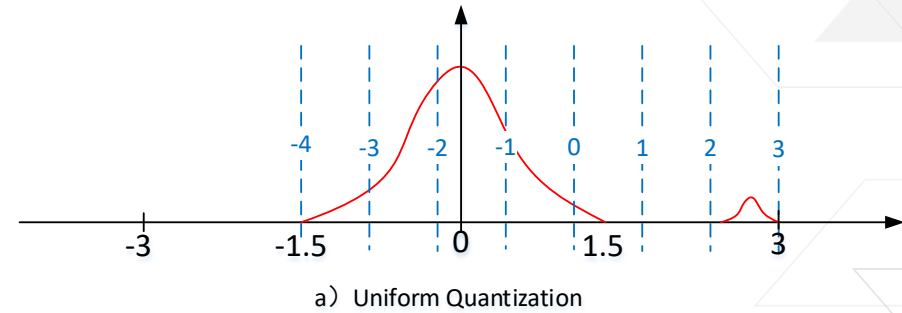
	Dynamic Range	Min Positive Value
FP32	$-3.4 \times 10^{38} \sim +3.4 \times 10^{38}$	1.4×10^{-45}
FP16	$-65504 \sim +65504$	5.96×10^{-8}
INT8	$-128 \sim +127$	1



Quantization Strategy

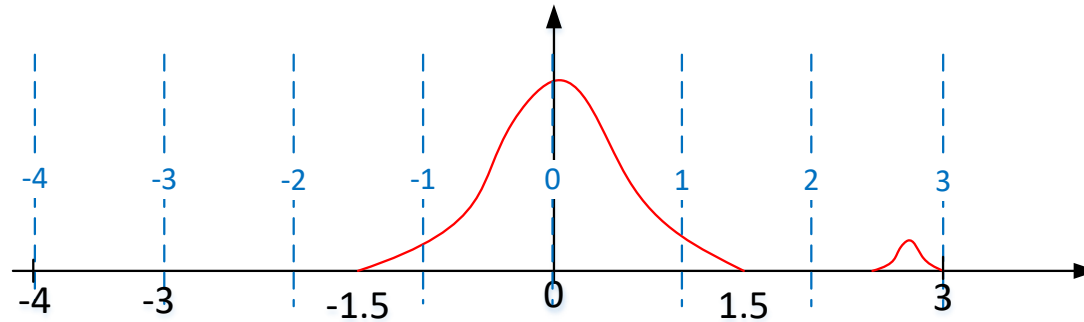
E.g.: FP32->INT3

- 8 discrete values (-4,-3,...,3)
- Optimization Target:
 - E.g. Non-overflow
 - E.g. Minimize L2 loss: $\sum (X-X')^2$
- Methods:
 - Uniform vs. Nonuniform
 - Symmetry vs. Asymmetry :
 - Zero-Shift(*TensorFlow*)
 - Float Scale vs. Int Scale



Our Strategy

- Quantization Strategy
 - Uniform Symmetric Quantization → 8Bit for DeePhi DPU
 - Scale = 2^N



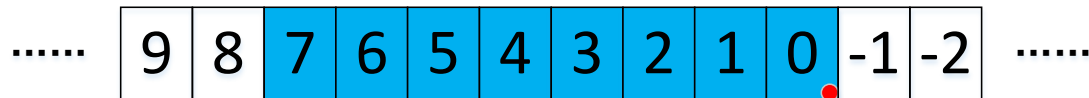
c) Uniform Symmetric Quantization

- Advantages
 - Simple hardware design
 - High efficiency: all fix-point calculation
 - Make use of redundancy of CNN models(especially with BatchNorm Layers)

Quantize, Dequantize, Quantize position

- Uniform Symmetric Quantization :
 - Scale: $2^{(-p)}$
 - Quantize:
 - $X_q = F(X_f) = \text{round}(X_f / \text{scale}) = \text{round}(2^{(p)} * X_f)$
 - Dequantize:
 - $X_d = F'(X_q) = X_q * \text{Scale} = 2^{(-p)} * X_q$
 - Range of dequantized value:
 - $[-128 * 2^{(-p)}, 127 * 2^{(-p)}]$
 - Quantize Pos:
 - $p = \log_2(1 / \text{AbsMax}(X_{\min}/128, X_{\max}/127))$

p	scale	min	max
-2	4	-512	508
-1	2	-256	254
0	1	-128	127
1	0.5	-64	63.5
2	0.25	-32	31.75
3	0.125	-16	15.875
4	0.0625	-8	7.9375



p: the decimal position of quantized value

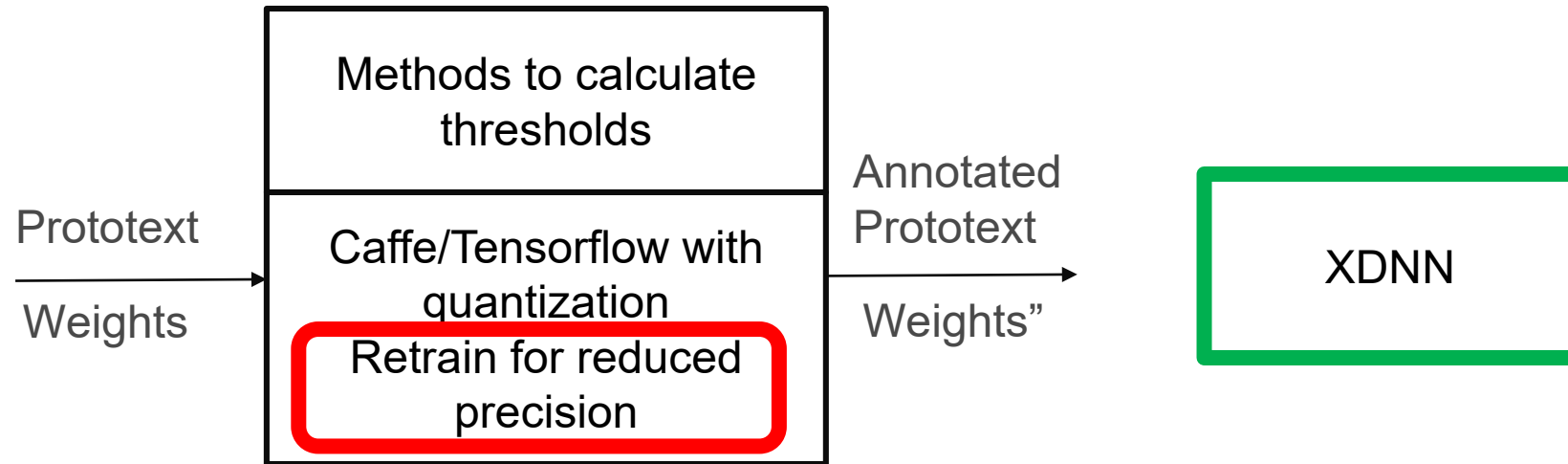
Example

- Example: Float→Int8
 - Float value: {-6.0625, -8.013, 4.438}
- Steps:
 - Float range: [-8.013~4.438]
 - Quantize pos: $p = 3$
 - Scale: $2^{(-3)} = 0.125$
 - Quantize: $x_q = \text{round}(x / \text{scale}) = \text{round}(x * 8)$
 - Dequantize: $x_d = x_q * \text{scale} = x_q * 0.125$

p	scale	min	max
-2	4	-512	508
-1	2	-256	254
0	1	-128	127
1	0.5	-64	63.5
2	0.25	-32	31.75
3	0.125	-16	15.875
4	0.0625	-8	7.9375

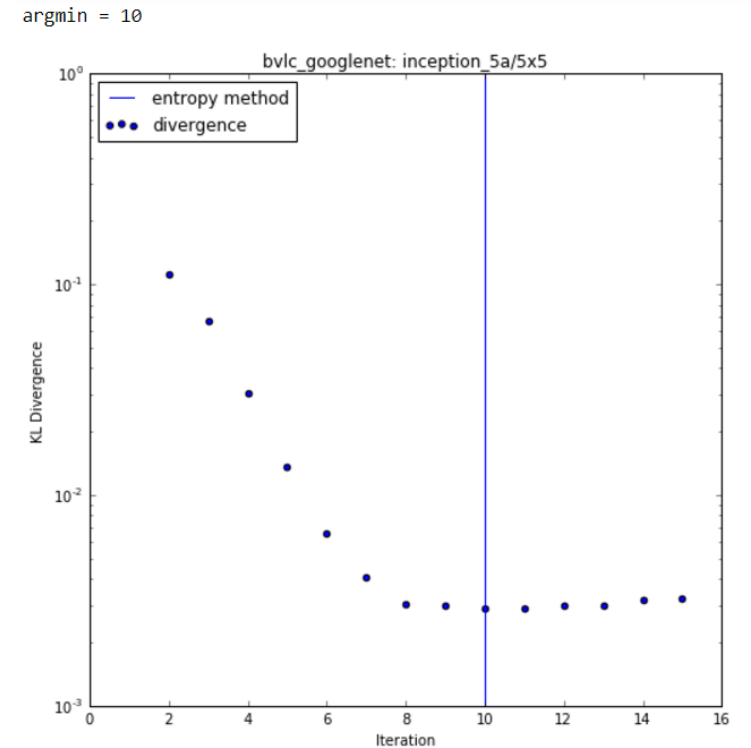
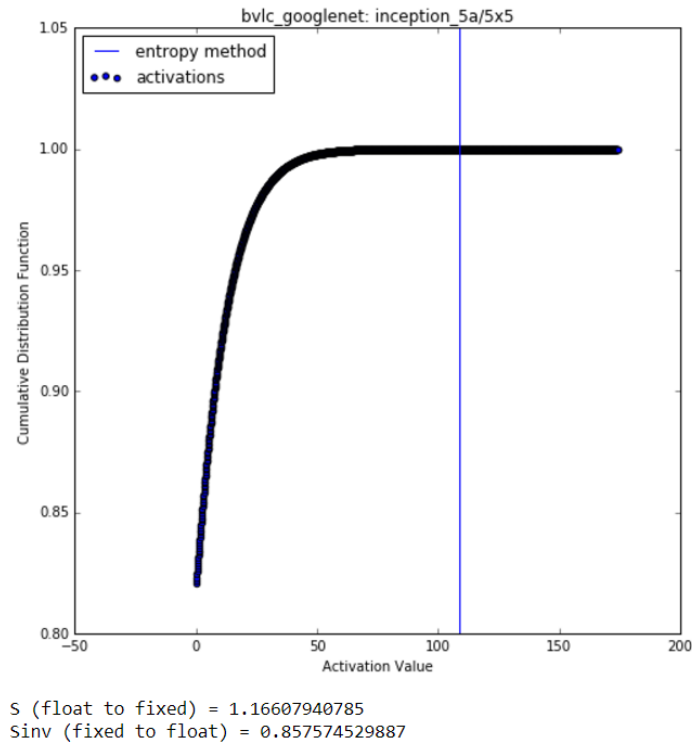
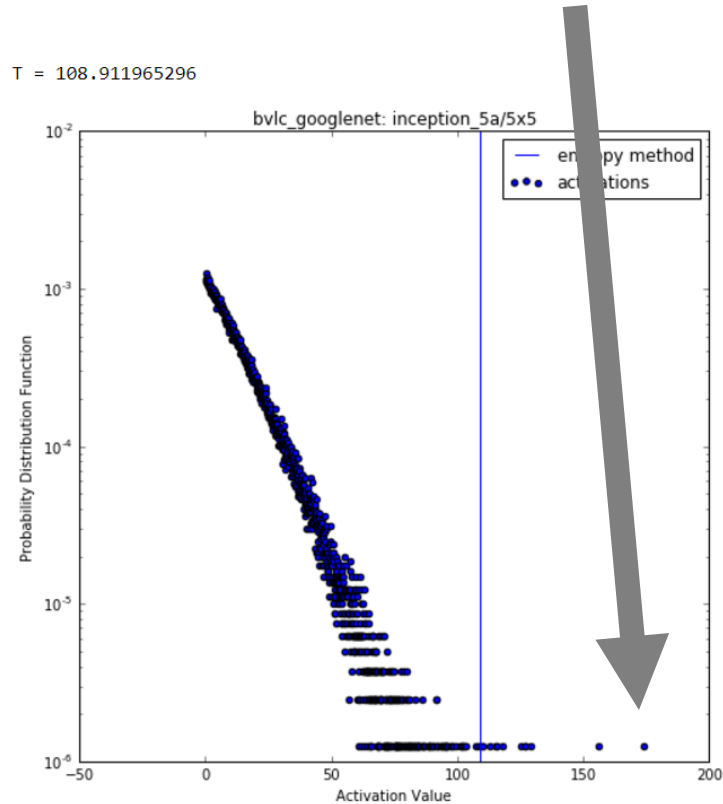
X	$X_q(p=3)$	X_d
-6.0625	-49	-6.125
-8.013	-64	-8
4.438	36	4.5

Overall Quantization Flow



Offline Quantization Procedure

- Walk the network graph, analyze distributions, and statistically determine thresholds for weights and activations using a calibration set
 - >> Where is the information in your deep learning models?



Example Results

Network	Dataset	Values	weights	Top1 %	diff	Top5 acc	diff
Resnet-50	ImageNet	Floating P.	Original	72.97	0.00	91.13	0.00
		8 bits	Quantized	70.67	-2.30	89.79	-1.34
		8 bits	Retrained	73.39	0.43	91.47	0.34
	Flowers102	Floating P.	Original	85.38	0.00	97.38	0.00
		8 bits	Quantized	81.57	-3.81	96.37	-1.00
		8 bits	Retrained	85.10	-0.28	97.06	-0.32
	Places365	Floating P.	Original	52.38	0.00	83.00	0.00
		8 bits	Quantized	50.99	-1.39	82.39	-0.61
		8 bits	Retrained	53.80	1.43	84.28	1.28

Retraining recovers accuracy, and uses initial weights from offline quantization
Retraining time is significantly lower than training from scratch.

xfDNN Runtime



xfDNN Runtime Engine

- > Lightweight and portable with no dependency on ML frameworks
 - >> Uses JSON exported by Compiler and Quantizer
 - >> Generate binary for xDNN
- > Extensive C++/Python API with simplified use model
- > Asynchronous XDNN execution
- > Streaming/pipeline/Dataflow support for large imageset
- > Multiple CNN models running on single FPGA
- > Multiple FPGA support
- > Support for both 16 and 8 bits
- > CPU accelerated functions



Simplify xdnn usage

```
import xdnn, xdnn_io

args = xdnn_io.processCommandLine()
xdnn.createHandle(args['xclbin'], "kernelSxdnn_0", args['xlnxlib'])

(weightsBlob, fcWeight, fcBias) = xdnn_io.loadWeights(args)

(fpgaInputs, batch_sz) = xdnn_io.prepareInput(args)
fpgaOutput = xdnn_io.prepareOutput(args['fpgaoutsz'], batch_sz)

xdnn.execute(args['netcfg'],
             weightsBlob, fpgaInputs, fpgaOutput,
             batch_sz, args['quantizecfg'])

fcOut = xdnn.computeFC(fcWeight, fcBias, fpgaOutput)

softmaxOut = xdnn.computeSoftmax(fcOut)

xdnn_io.printClassification(softmaxOut, args);

xdnn.closeHandle()
```

Blocking

```
# exec_async() enqueues FPGA task and returns immediately
xdnn.exec_async(args['netcfg'],
               weightsBlob, fpgaInputs, fpgaOutput,
               batch_sz, args['quantizecfg'], args['PE'])

# get_result() blocks for FPGA result
xdnn.get_result(args['PE'])
```

Non-Blocking

```
numFPGAs = 8
args = xdnn_io.processCommandLine()
xdnn.createHandle(args['xclbin'], "kernelSxdnn_0", args['xlnxlib'], numFPGAs)

(weightsBlob, fcWeight, fcBias) = xdnn_io.loadWeights(args)
```

8-FPGA

Streaming/pipeline/dataflow

```
def main():
    qPrep = Queue(maxsize=1)
    qFpga = Queue(maxsize=1)

    prepProc = Process(target=prep_process, args=(qPrep,))
    xdnnProc = Process(target=xdnn_process, args=(qPrep, qFpga))
    postProc = Process(target=post_process, args=(qFpga,))

    prepProc.start()
    xdnnProc.start()
    postProc.start()
```

```
def xdnn_process(qFrom, qTo):
    xdnn.createHandle(g_xclbin, "kernelSxdnn_0", g_xdnnLib)
    weightsBlob = loadWeightsBiasQuant()

    fpgaOutput = None
    while True:
        (inputs, inputImageFiles) = qFrom.get()
        fpgaInputs = prepareFpgaInputs(inputs)
        if not fpgaOutput:
            fpgaOutput = prepareOutput(g_batchSize)

        xdnn.execute(g_netFile,
                    weightsBlob, fpgaInputs, fpgaOutput,
                    g_batchSize, g_fpgaCfgFile)

        np_out = np.frombuffer(fpgaOutput, np.float32)
        qTo.put((np_out, inputImageFiles))
```

```
def prep_process(q):
    ret = xdnn.createManager(g_xdnnLib)

    cInputBuffer = None
    cFpgaInputBuffer = None
    while True:
        (inputs, inputImageFiles) = prepareImages()

        fpgaInputs = xdnn.quantizeInputs(g_firstFpgaLayerName,
                                         inputs, cInputBuffer, cFpgaInputBuffer, g_fpgaCfgFile)

        q.put((fpgaInputs, inputImageFiles))
```

```
def post_process(qFpga):
    (fcWeight, fcBias) = loadFCWeightsBias()

    while True:
        (fpgaOutput, inputImageFiles) = qFpga.get()

        fcOutput = xdnn.computeFC(fcWeight, fcBias, fpgaOutput)
        smaxOutput = xdnn.computeSoftmax(fcOutput)

        printClassification(smaxOutput, inputImageFiles, g_labelFile)
```

Multiple CNN models running on single FPGA

```
args = xdn_io.processCommandLine()
xdnn.createHandle(args['xclbin'], "kernelSxdnn_0", args['xlnxlib'])

# setup/pre-process
for netCfg in args['fpgacfg']:
    netID = netCfg['name']
    PE = [int(x) for x in netCfg['PE'].split()]
    PEs[netID] = PE

    (weightsBlobs[netID], fcWeights[netID], fcBiases[netID]) = xdn_io.loadWeights(netCfg)
    (fpgaInputs[netID], batchSize) = xdn_io.prepareInput(netCfg)
    batchSize[netID] = batchSize
    fpgaOutputs[netID] = xdn_io.prepareOutput(int(netCfg['fpgaoutsz']), batchSize)
    netFiles[netID] = netCfg['netcfg']

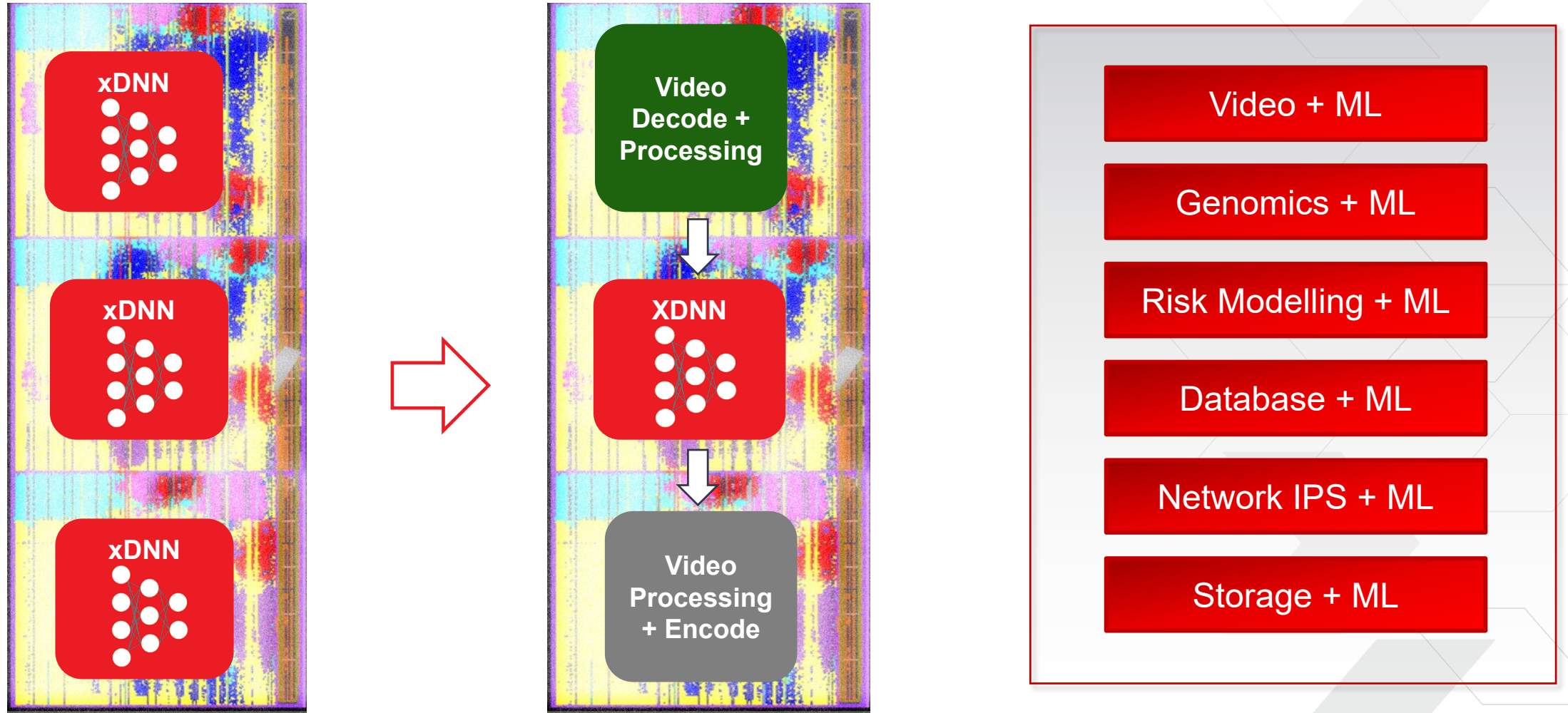
# exec different networks on different PEs (async)
for netCfg in args['fpgacfg']:
    netID = netCfg['name']
    xdn.exec_async(netFiles[netID], weightsBlobs[netID], fpgaInputs[netID],
                  fpgaOutputs[netID], batchSize[netID], 1, netCfg['quantizecfg'], PEs[netID])

# collect results
for netID, pe in PEs.iteritems():
    xdn.get_result(pe)

# post-process
for netCfg in args['fpgacfg']:
    netID = netCfg['name']
    fcOut = xdn.computeFC(fcWeights[netID], fcBiases[netID], fpgaOutputs[netID])
    softmaxOut = xdn.computeSoftmax(fcOut)
    xdn_io.printClassification(softmaxOut, netCfg);
```

```
fpgacfg:
[
  {
    "name": "googlenet_1",
    "net": "googlenet",
    "datadir": "data",
    "netcfg": "googlenet.fpgaaddr.64.txt",
    "PE": "1",
    "firstFpgaLayerName": "conv1/7x7_s2",
    "fpgaoutsz": "1024"
  },
  {
    "name": "resnet_2",
    "net": "resnet",
    "datadir": "data_resnet",
    "netcfg": "resnet.fpgaaddr.64.txt",
    "PE": "2",
    "firstFpgaLayerName": "conv1/7x7_s2",
    "fpgaoutsz": "2048"
  }
]
```

Deep Learning Infusion



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

Adaptable.
Intelligent.

