

Compute-Efficient Neural-Network Acceleration

Ephrem Wu, Xiaoqian Zhang, David Berman, Inkeun Cho, John Thendean

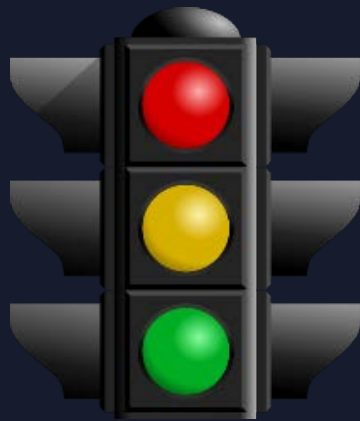
FPGA 2019

2/24/2019

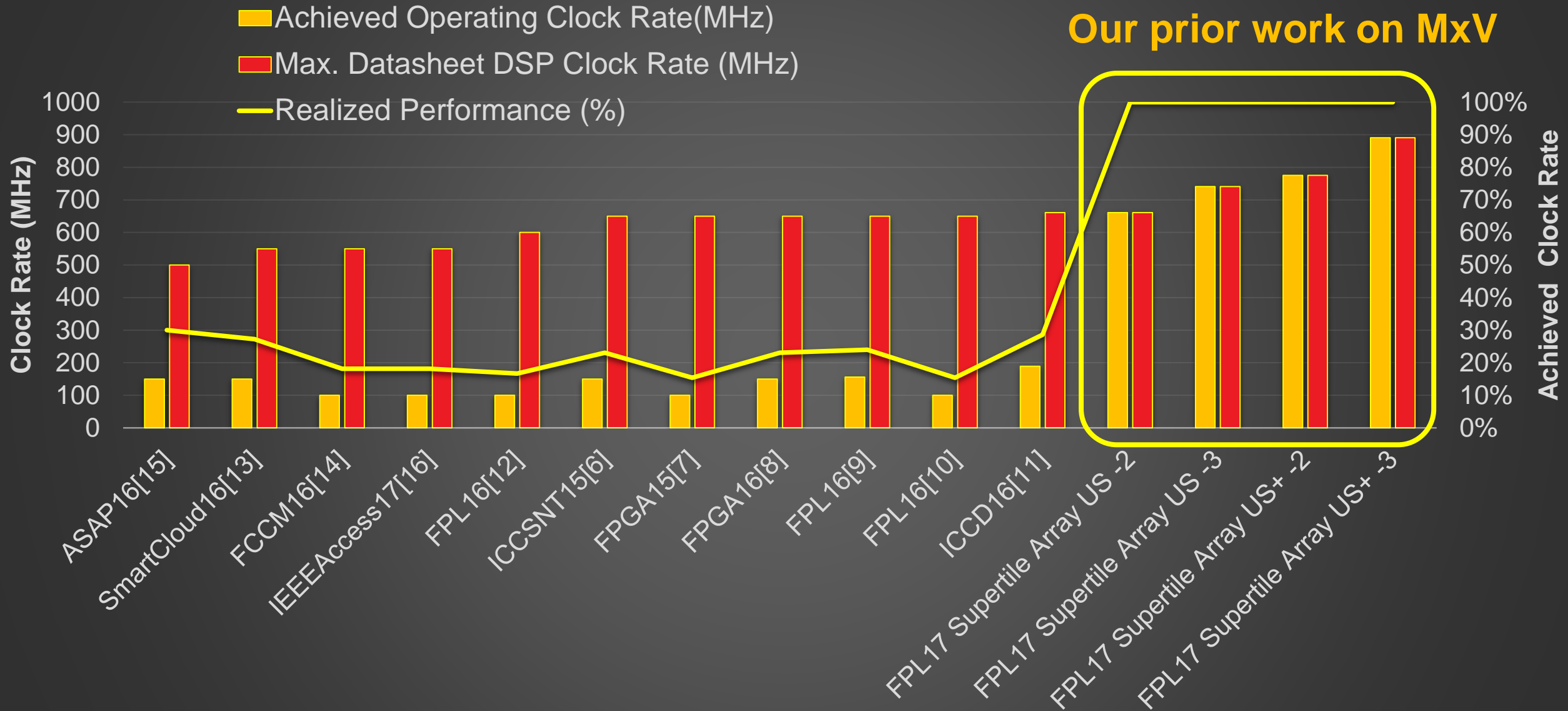


Goal

Build a compute-efficient and reconfigurable neural-net inference accelerator



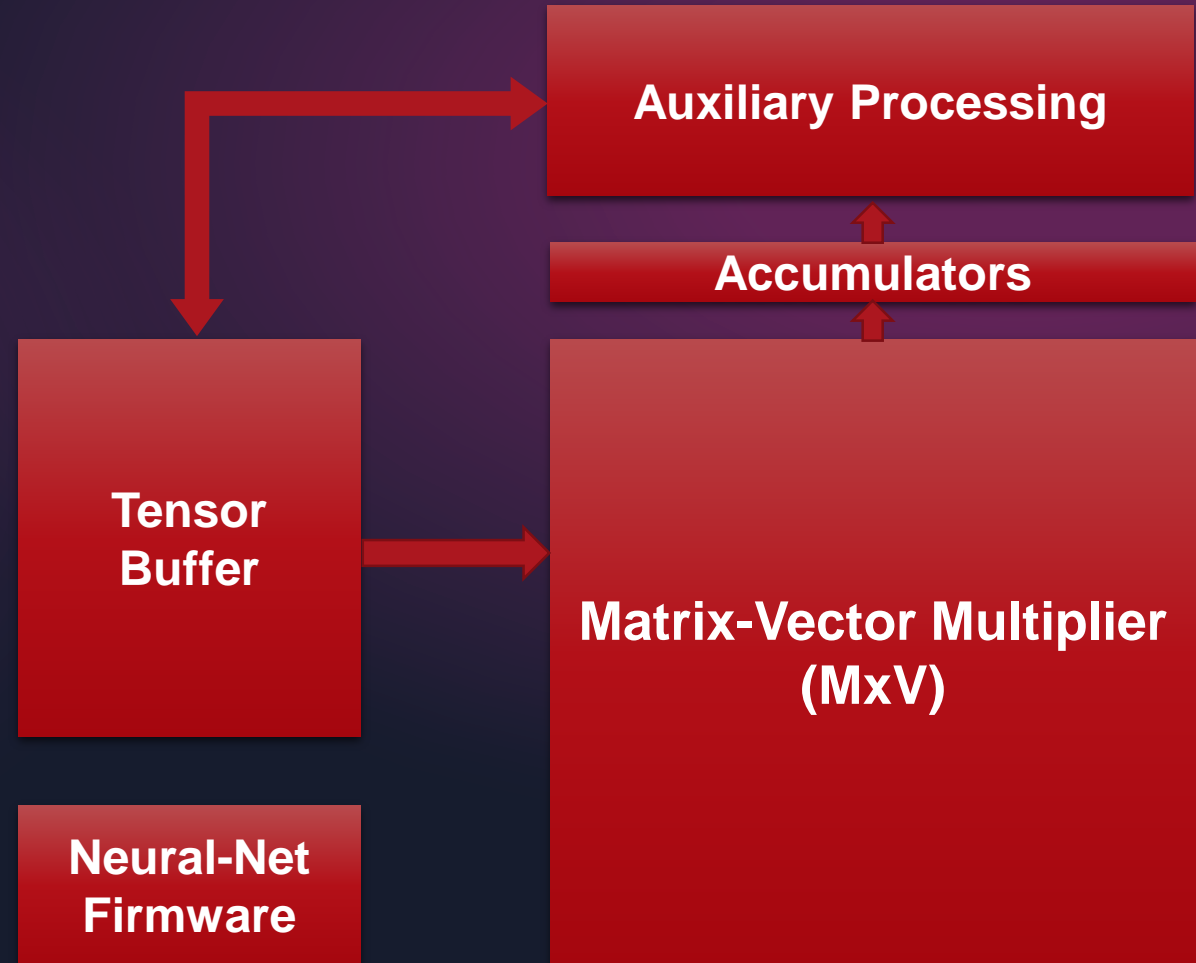
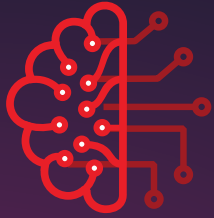
We Don't Have to Leave DSP Performance Behind



Key Design Principles

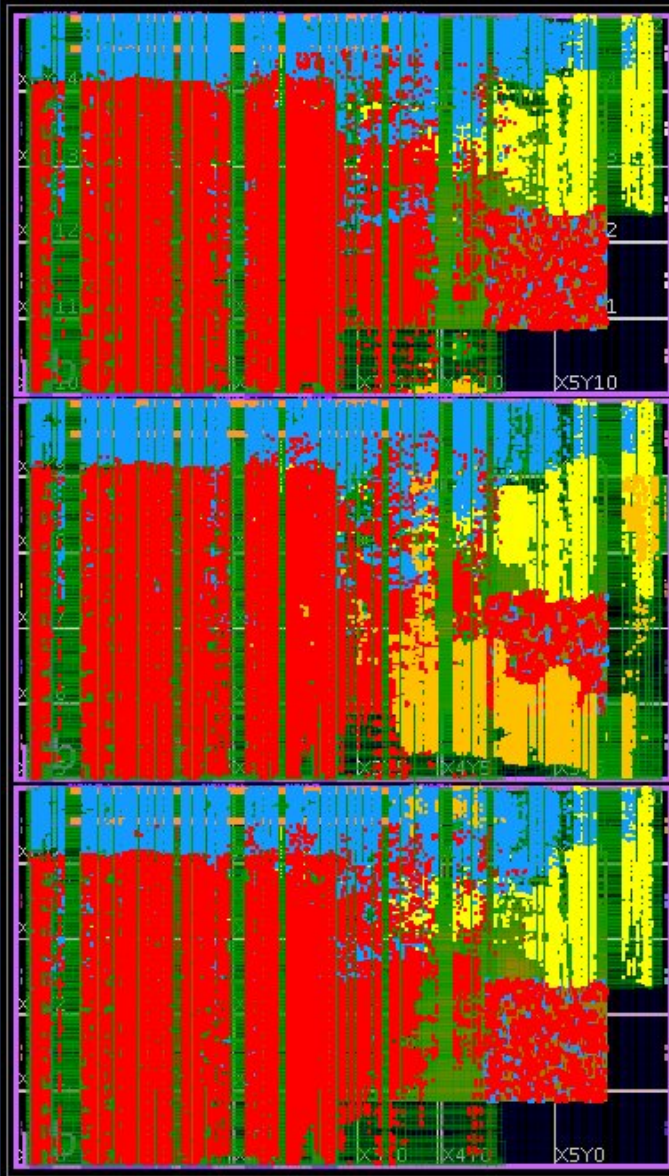
1. **Make compute fast.**
2. **Keep compute busy.**
3. **Simplify implementation.**

A Domain-Specific Architecture for Neural Networks



Domain-specific architectures called for by Hennessey and Patterson. See "A New Golden Age for Computer Architecture," [Communications of the ACM](#), February 2019, Vol. 62 No. 2, Pages 48-60.

Compute-Efficient Neural Processing



VU9P Layout

> **Case Study: GoogLeNet v1 Inference**

- >> 3 parallel GoogLeNets with their own weights
- >> Aggregate 3046 images/sec, 3.3 ms latency

> **Platform: VCU1525 board with VU9P-2 FPGA**

> **Compute**

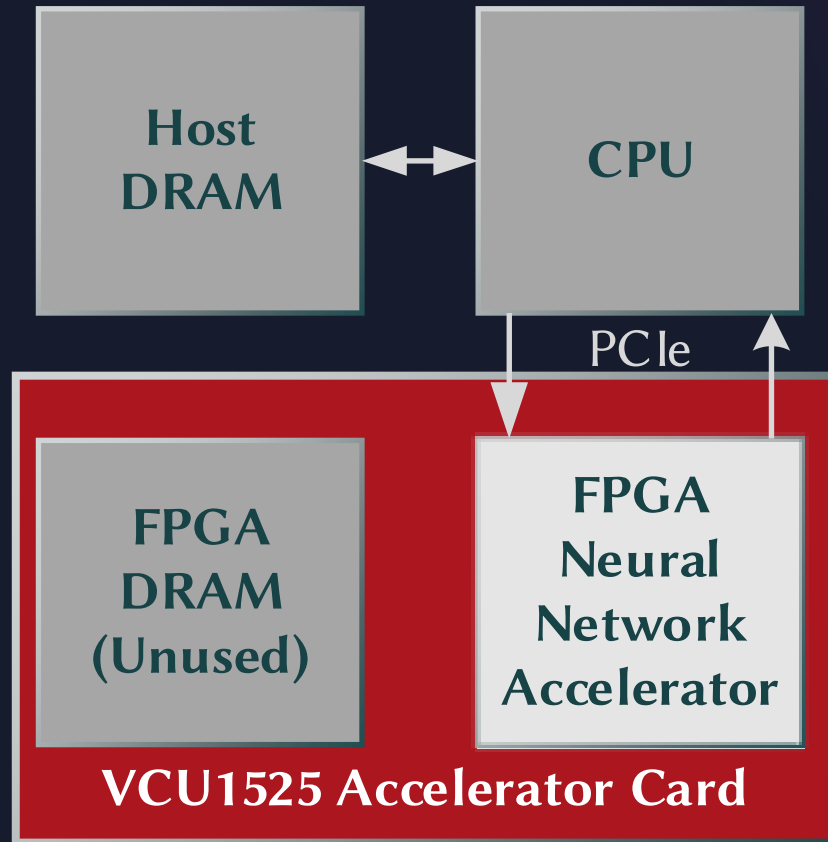
- >> DSP supertile arrays running at 720 MHz
- >> Consumes only 56% DSP48 tiles
- >> DSP cycles 95% utilized
- >> Per-tensor block floating-point, 8-/16-bit significands

> **Memory**

- >> No external DRAM on accelerator card used
- >> All tensors stored in UltraRAM & BRAM
- >> 1/2 DSP clock rate

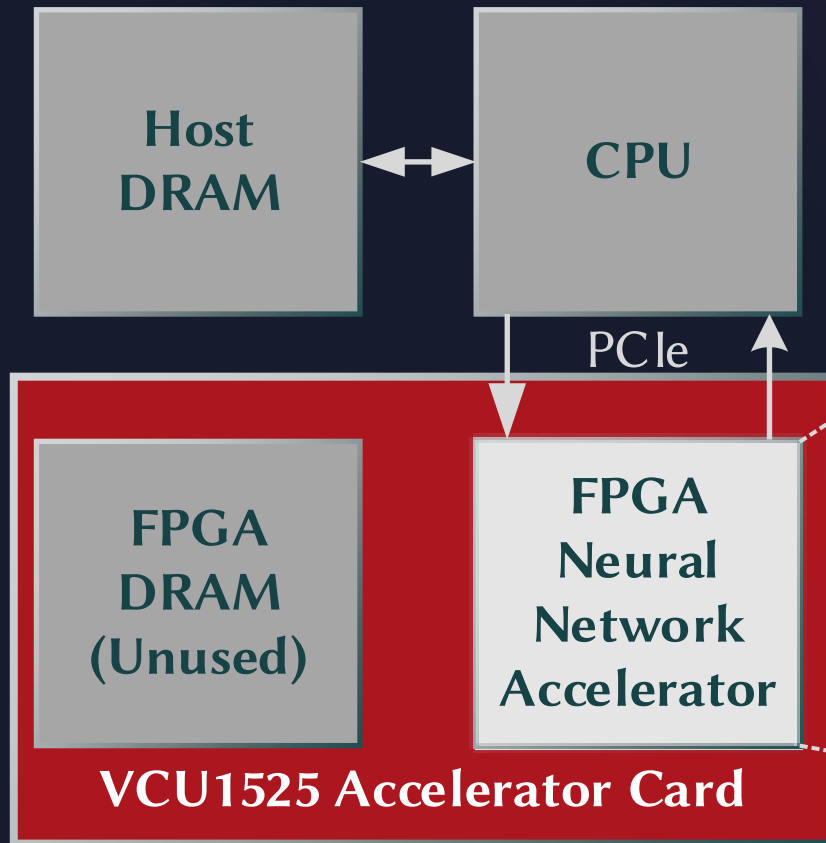
System and Accelerator Block Diagrams

System

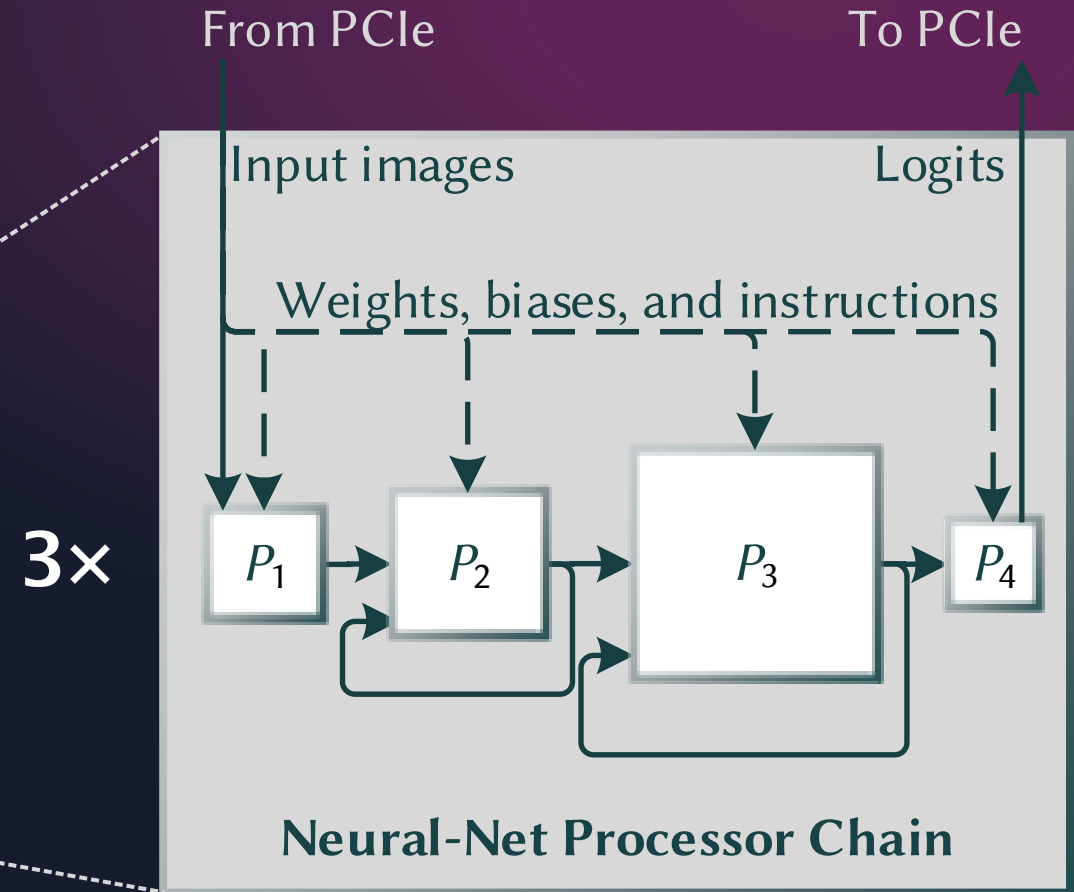


System and Accelerator Block Diagrams

System



Accelerator



Apples-to-Apples Accelerator Comparison

- > Throughput per device doesn't tell the whole story
- > Differences: FPGAs, total compute, data movement, numerical formats, and clock rates
- > Normalize results
 - >> Achieved Clock Rate: Actual DSP clock compared to datasheet F_{\max}
 - >> Compute Efficiency: % available DSP cycles that do useful work

Make compute fast

Keep compute busy

Overall Efficiency = % DSP F_{\max} × Compute Efficiency

GoogLeNet Inference Comparison

FPGA Implementation	Realized Clock Rate (% of Peak)	Compute Efficiency	Overall Efficiency
Ngo, 2016	15.4%	15.0%	2.3%
Venieris, 2017	19.2%	66.6%	12.8%
Shen, 2017	26.1%	93.8%	24.5%
Huang, 2018	30.8%	81.4%	25.1%
Gokhale, 2017	38.4%	91.0%	34.9%
This work	92.9%	95.5%	88.7%

The higher the clock rate, the harder it is to reach high compute efficiency

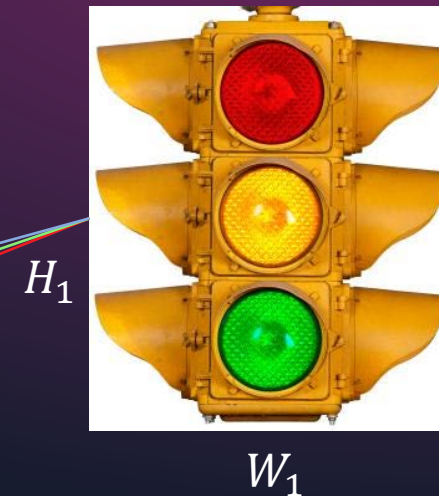
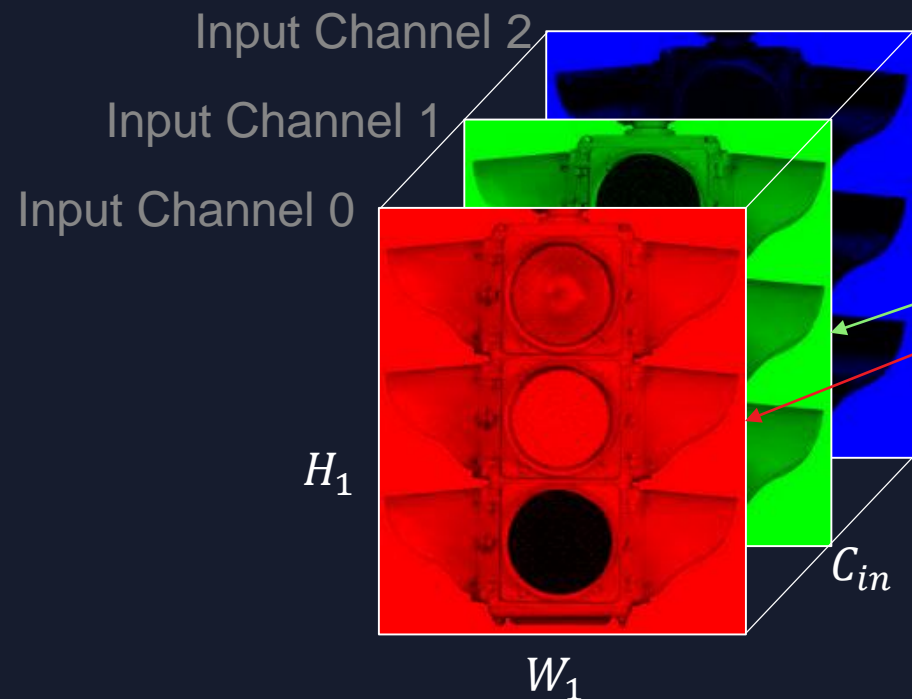
Bang for the buck from each DSP resource

Convolution: From Tensors View to Matrix View



2D Convolution: Tensor View

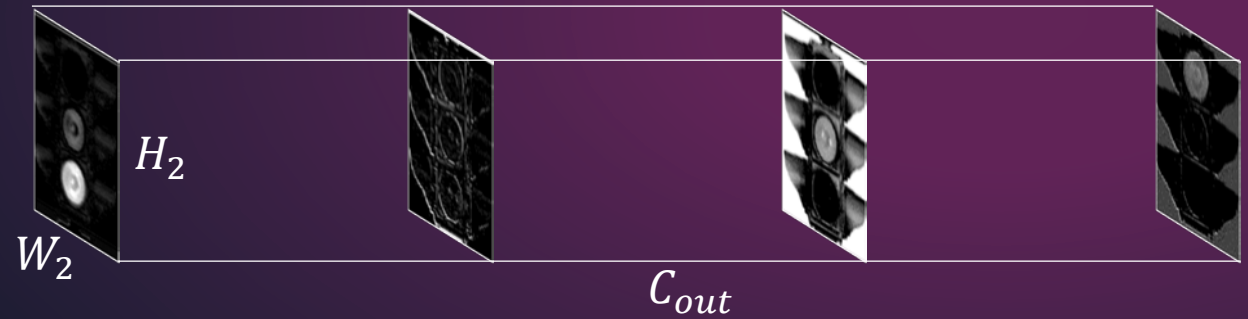
Input Tensor $\mathcal{X} \in \mathbb{R}^{H_1 \times W_1 \times C_{in}}$



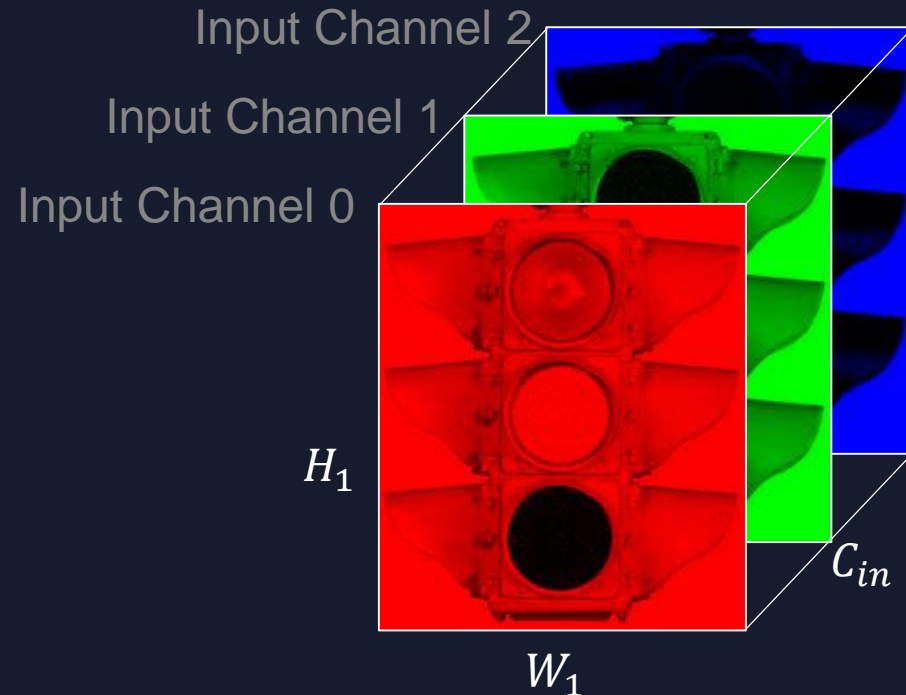
2D Convolution: Tensor View

Output Tensor $\mathcal{Y} \in \mathbb{R}^{H_2 \times W_2 \times C_{out}}$

Output channel 0 Output channel 1 Output channel 2 Output channel 3



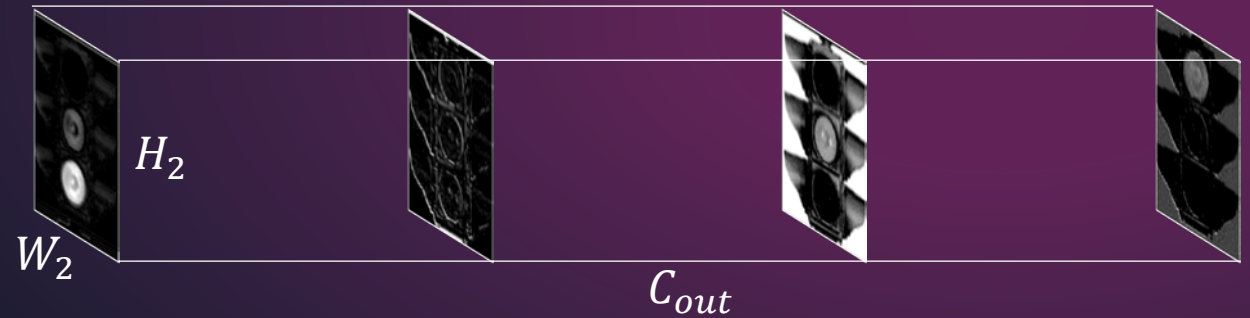
Input Tensor $\mathcal{X} \in \mathbb{R}^{H_1 \times W_1 \times C_{in}}$



2D Convolution: Tensor View

Output Tensor $\mathcal{Y} \in \mathbb{R}^{H_2 \times W_2 \times C_{out}}$

Output channel 0 Output channel 1 Output channel 2 Output channel 3



Input Tensor $\mathcal{X} \in \mathbb{R}^{H_1 \times W_1 \times C_{in}}$

Input Channel 2

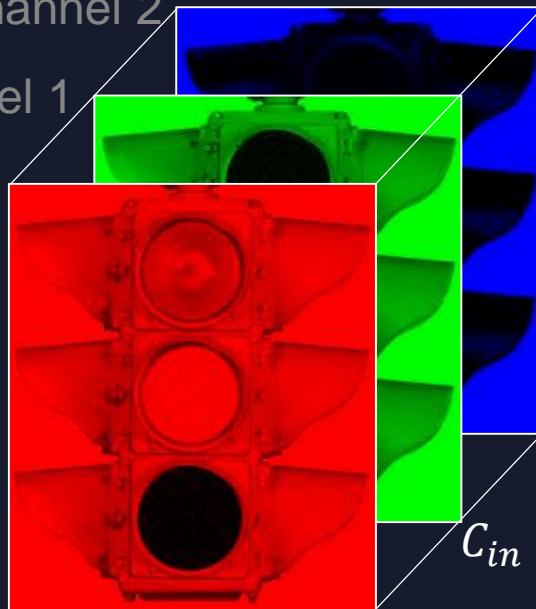
Input Channel 1

Input Channel 0

H_1

C_{in}

W_1



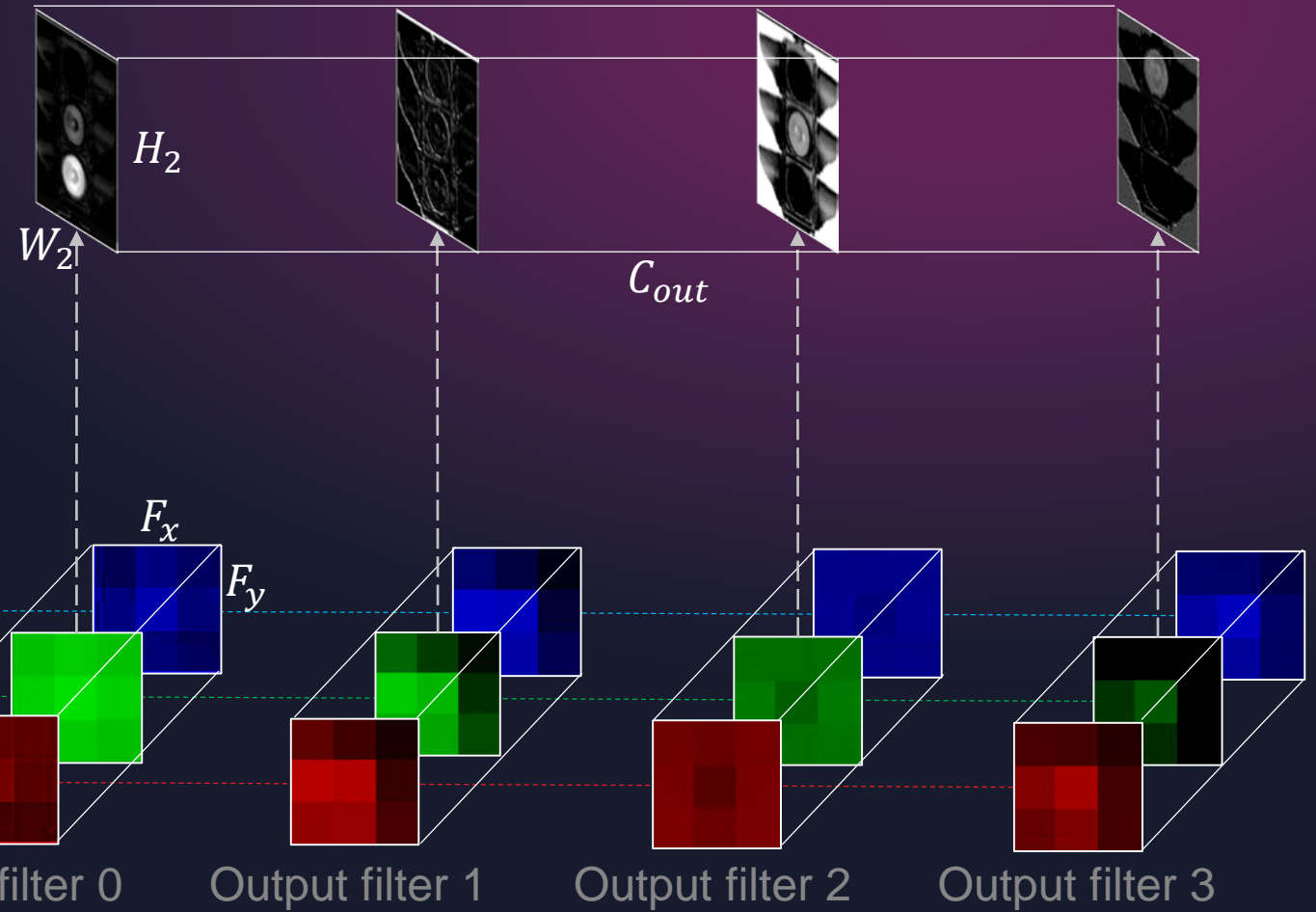
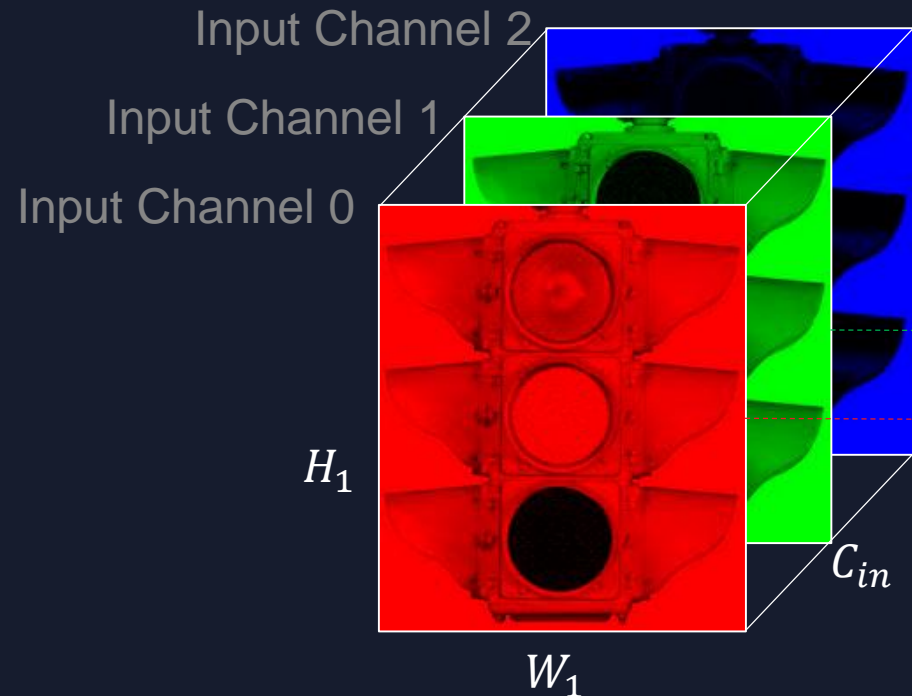
- > Although input and output tensors have three axes, convolution is 2D, not 3D.
- > There's a 2D convolution filter mask for every input-output channel pair ($3 \times 4 = 12$ in this example).

2D Convolution: Tensor View

Output Tensor $\mathcal{Y} \in \mathbb{R}^{H_2 \times W_2 \times C_{out}}$

Output channel 0 Output channel 1 Output channel 2 Output channel 3

Input Tensor $\mathcal{X} \in \mathbb{R}^{H_1 \times W_1 \times C_{in}}$



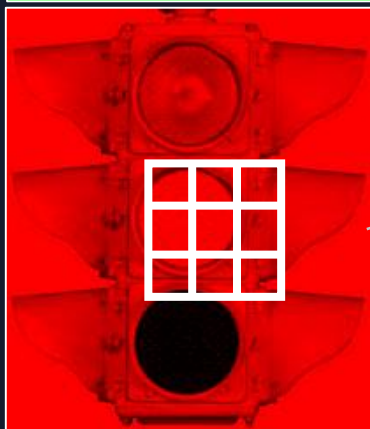
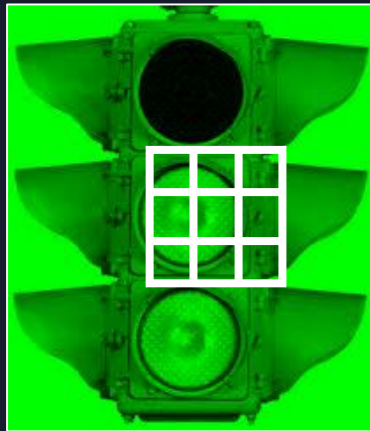
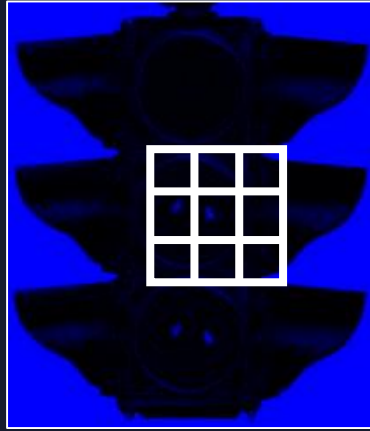
Filter Weight Tensor $\mathcal{W} \in \mathbb{R}^{F_y \times F_x \times C_{in} \times C_{out}}$

2D Convolution: From Tensor to Flat View

Input Image



Input Channels



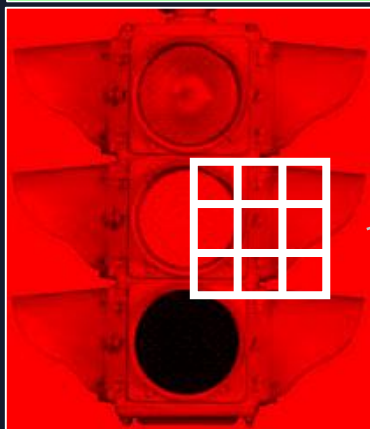
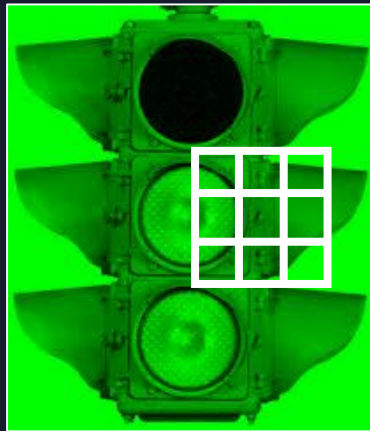
Parallel convolution sliding windows, one per input channel

2D Convolution: Flat View

Input Image



Input Channels



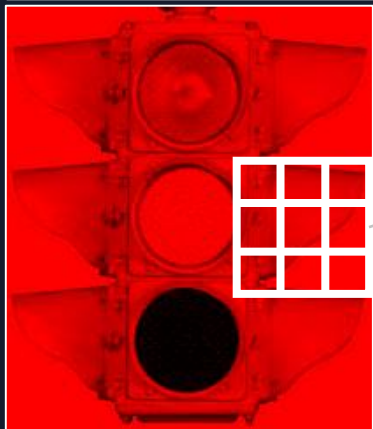
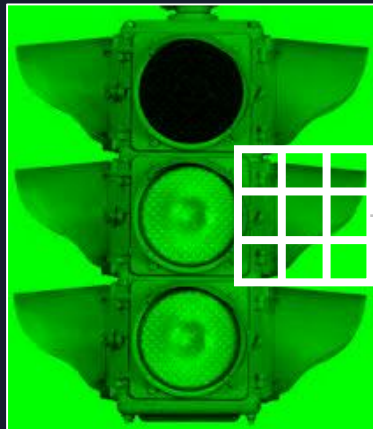
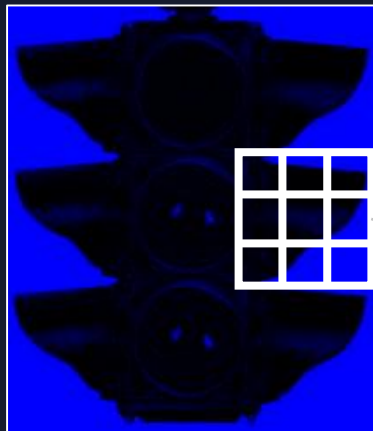
Parallel convolution sliding windows, one per input channel

2D Convolution: Flat View

Input Image



Input Channels



Parallel convolution sliding windows, one per input channel

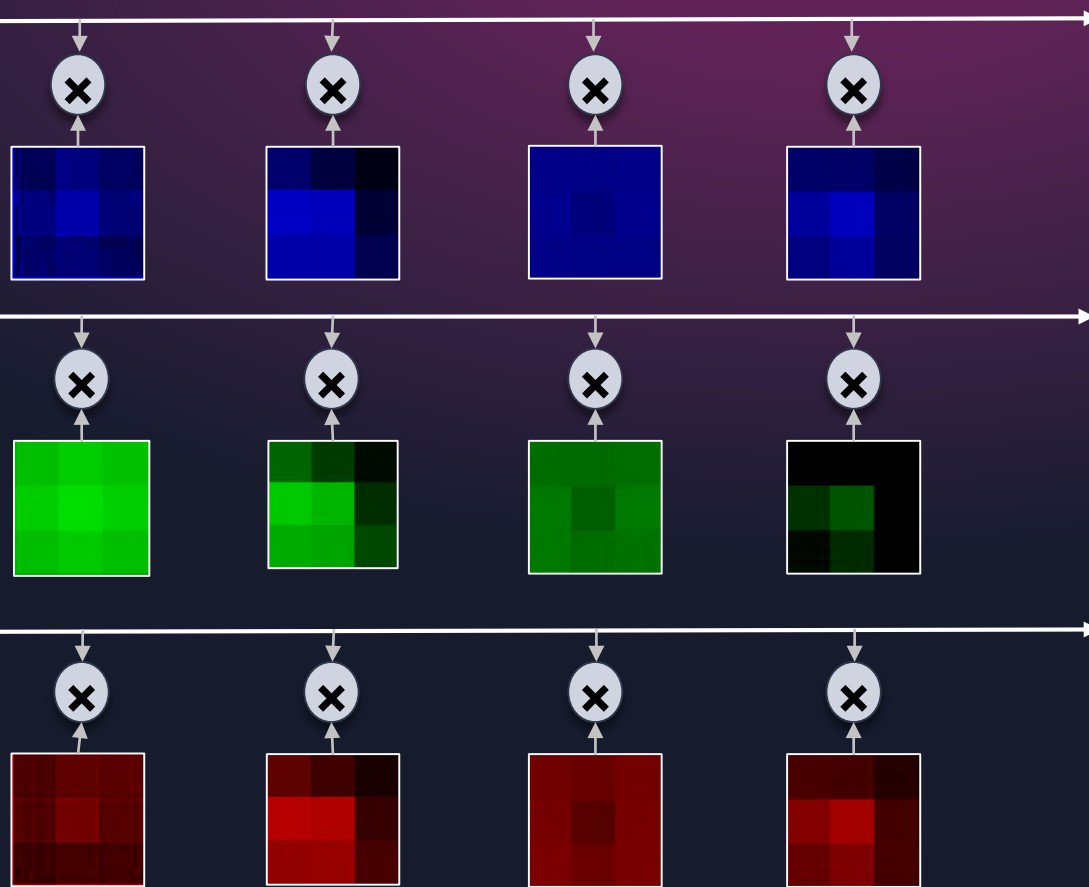
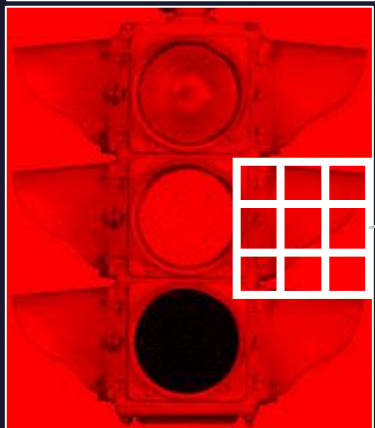
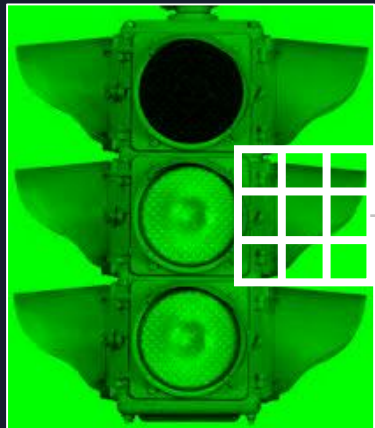
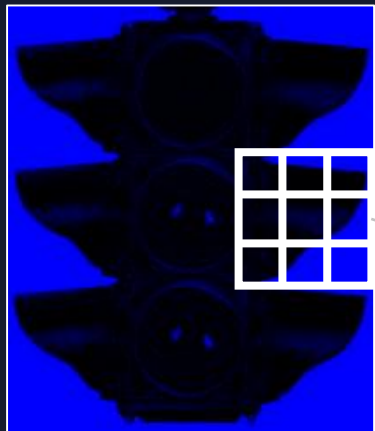
2D Convolution: Flat Silicon View

- Input channel broadcast
- Per-channel element-wise multiplication with filter weights

Input Image



Input Channels



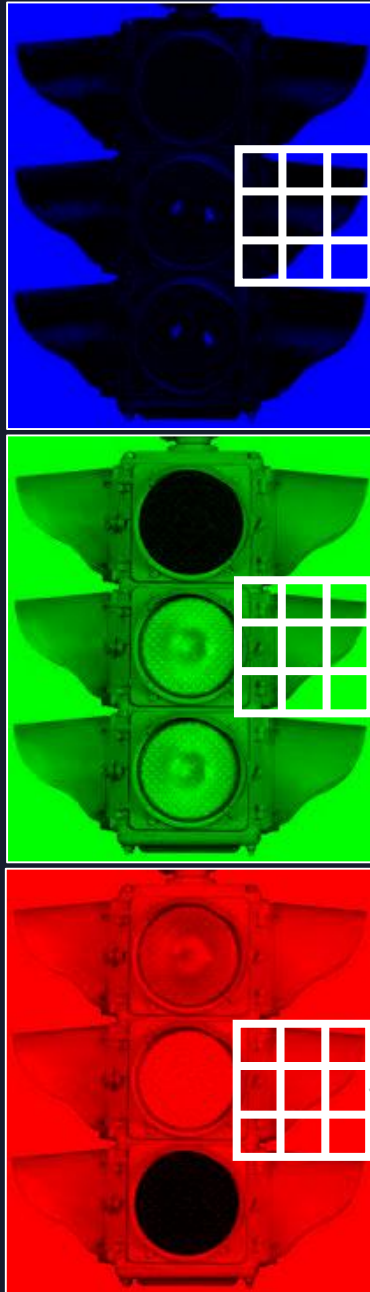
2D Convolution: Flat Silicon View

Reduction per output channel

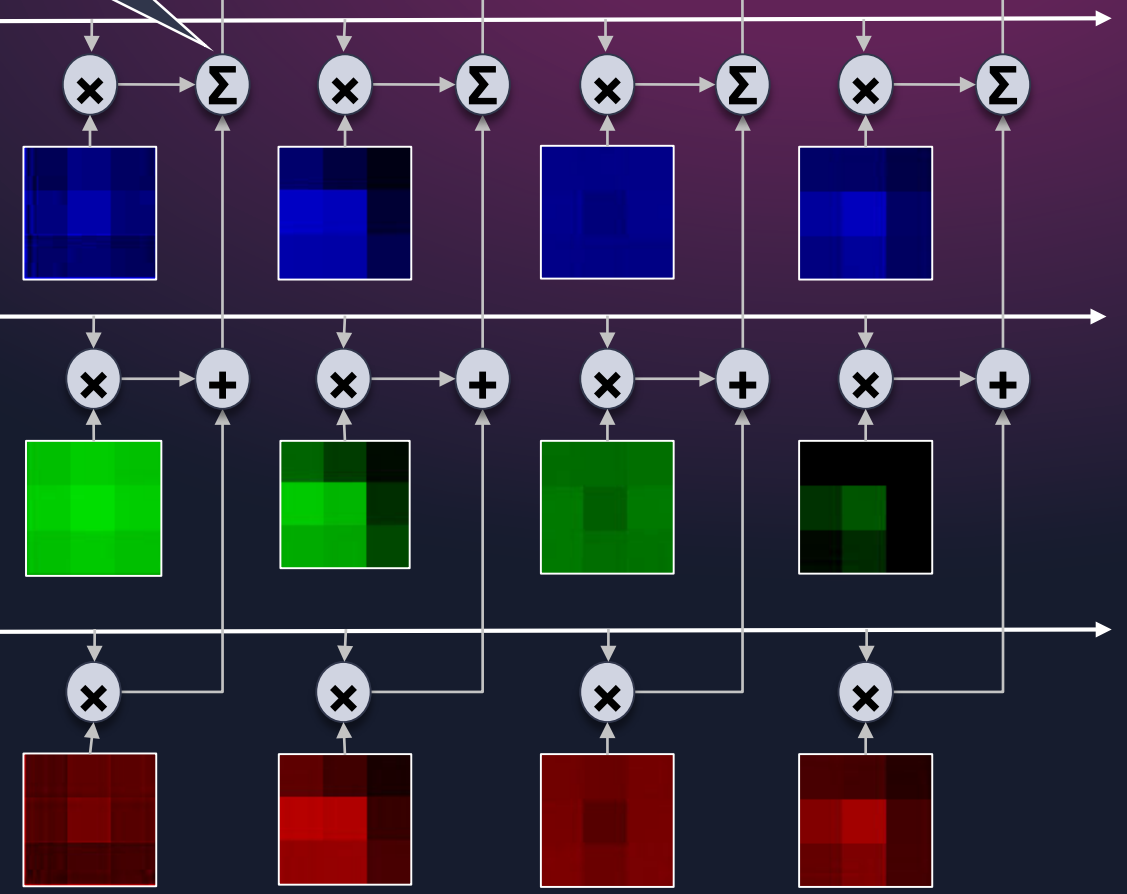
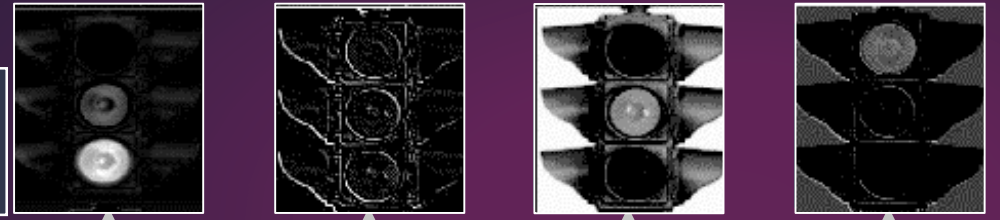
Input Image



Input Channels



Accumulates
 $3 \times 3 \times 3 = 27$ products

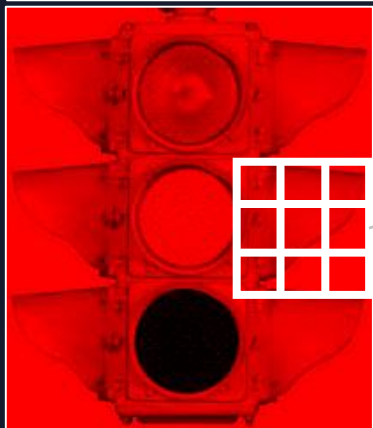
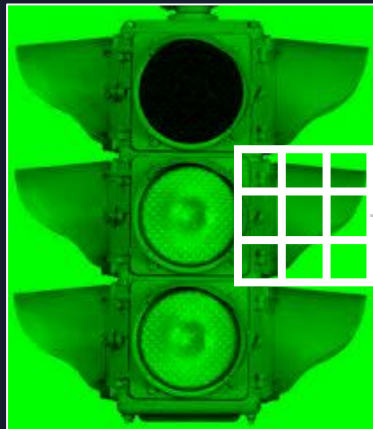
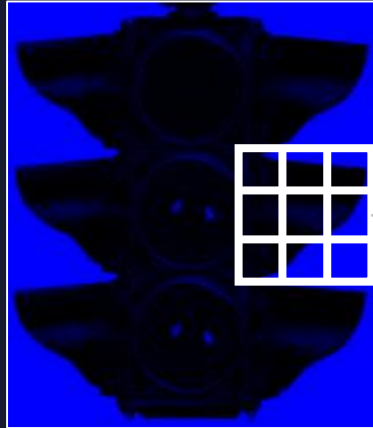


Interpreting Weights (Sort of)

Input Image

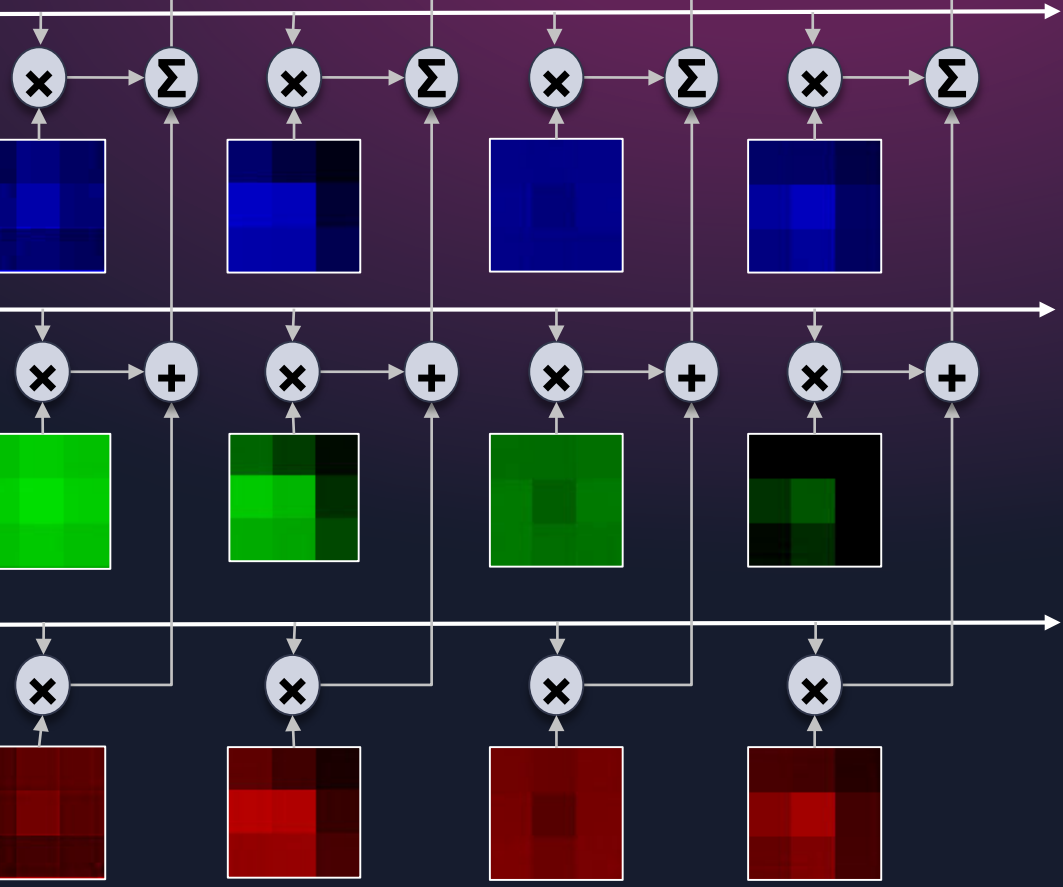
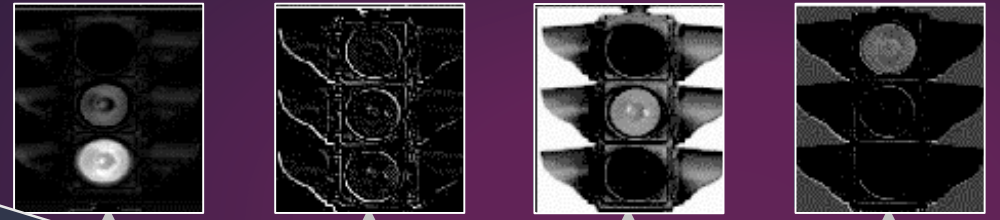


Input Channels



Looking for edges

Looking for green patches



Computing Convolution: Tensor Matricization and MxV



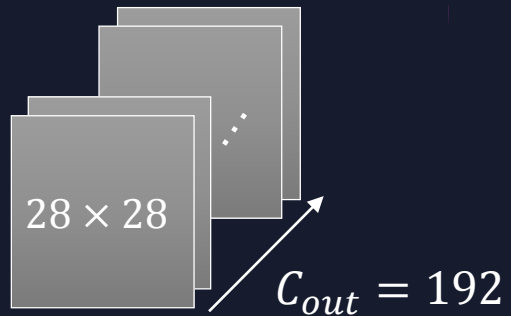
Computing Convolution

- > For each element in some output channel
 - >> Input vector: Flattened input patches
 - >> Weight vector: Flattened filter weights for this output channel
 - >> Compute dot product between input vector and weight vector
- > Input vector broadcast to multiple output channels
 - Hardware matrix-vector multiplier ($M \times V$)

Reshaping Tensors to Matrices

3-Mode Output Tensor

$$y \in \mathbb{R}^{H_2 \times W_2 \times C_{out}}$$



Elements per Output Channel

$$N = 28 \times 28 = 784$$

Output channels

$$M = 192$$

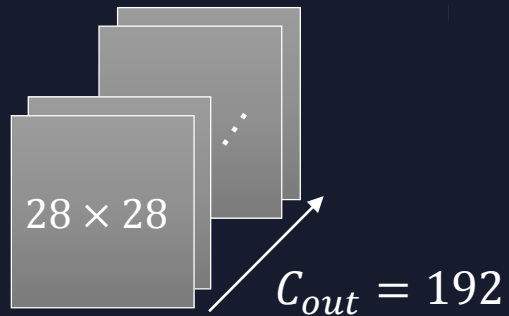
Output Matrix

$$Y \in \mathbb{R}^{C_{out} \times H_2 W_2}$$

Reshaping Tensors to Matrices

3-Mode Output Tensor

$$Y \in \mathbb{R}^{H_2 \times W_2 \times C_{out}}$$



Elements per Output Channel

$$N = 28 \times 28 = 784$$

Output channels

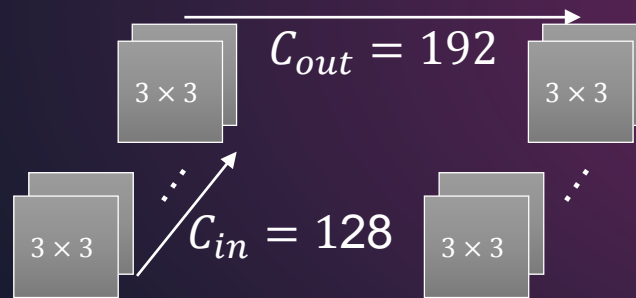
$$M = 192$$

Output Matrix

$$Y \in \mathbb{R}^{C_{out} \times H_2 W_2}$$

4-Mode Weight Tensor

$$W \in \mathbb{R}^{F_y \times F_x \times C_{in} \times C_{out}}$$



Input Channels \times Weights per Filter

$$K = 3 \times 3 \times 128 = 1152$$

$$M = 192$$

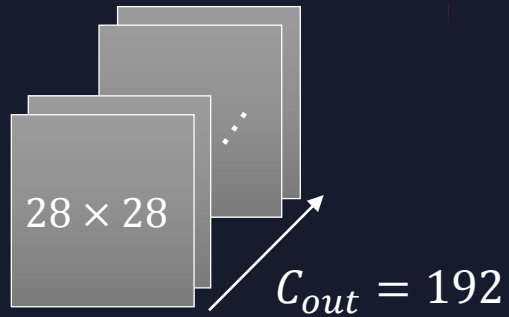
Weight Matrix

$$W \in \mathbb{R}^{C_{out} \times F_y F_x C_{in}}$$

Reshaping Tensors to Matrices

3-Mode Output Tensor

$$Y \in \mathbb{R}^{H_2 \times W_2 \times C_{out}}$$



Elements per Output Channel

$$N = 28 \times 28 = 784$$

Output channels

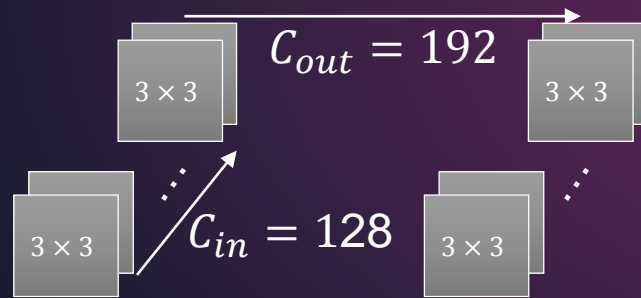
$$M = 192$$

Output Matrix

$$Y \in \mathbb{R}^{C_{out} \times H_2 W_2}$$

4-Mode Weight Tensor

$$W \in \mathbb{R}^{F_y \times F_x \times C_{in} \times C_{out}}$$



Input Channels \times Weights per Filter

$$K = 3 \times 3 \times 128 = 1152$$

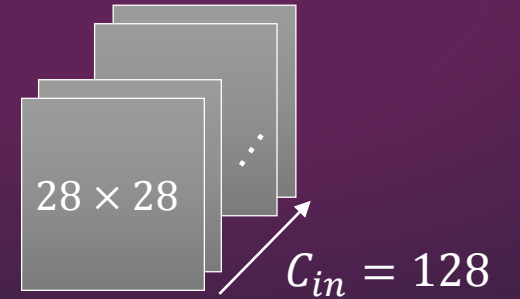
$$M = 192$$

Weight Matrix

$$W \in \mathbb{R}^{C_{out} \times F_y F_x C_{in}}$$

3-Mode Input Tensor

$$X \in \mathbb{R}^{H_1 \times W_1 \times C_{in}}$$



$$N = 28 \times 28 = 784$$

$$K = 1152$$

Input Matrix

$$X \in \mathbb{R}^{F_y F_x C_{in} \times H_2 W_2}$$

Example Weight Tensor Matricization Using Input Channel Interleaving

Input channel Input channel Input channel

R

1	4	7
2	5	8
3	6	9

G

1	4	7
2	5	8
3	6	9

B

1	4	7
2	5	8
3	6	9

Output channel j

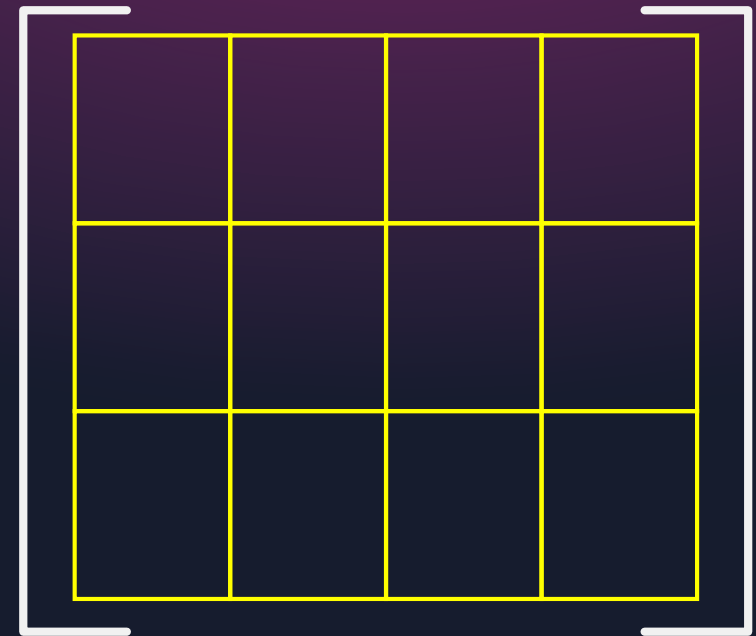
Ordering:
Weight column, weight row,
input channel

Input Channels \times Weights per Filter



Distributing Compute Over Hardware MxV Array

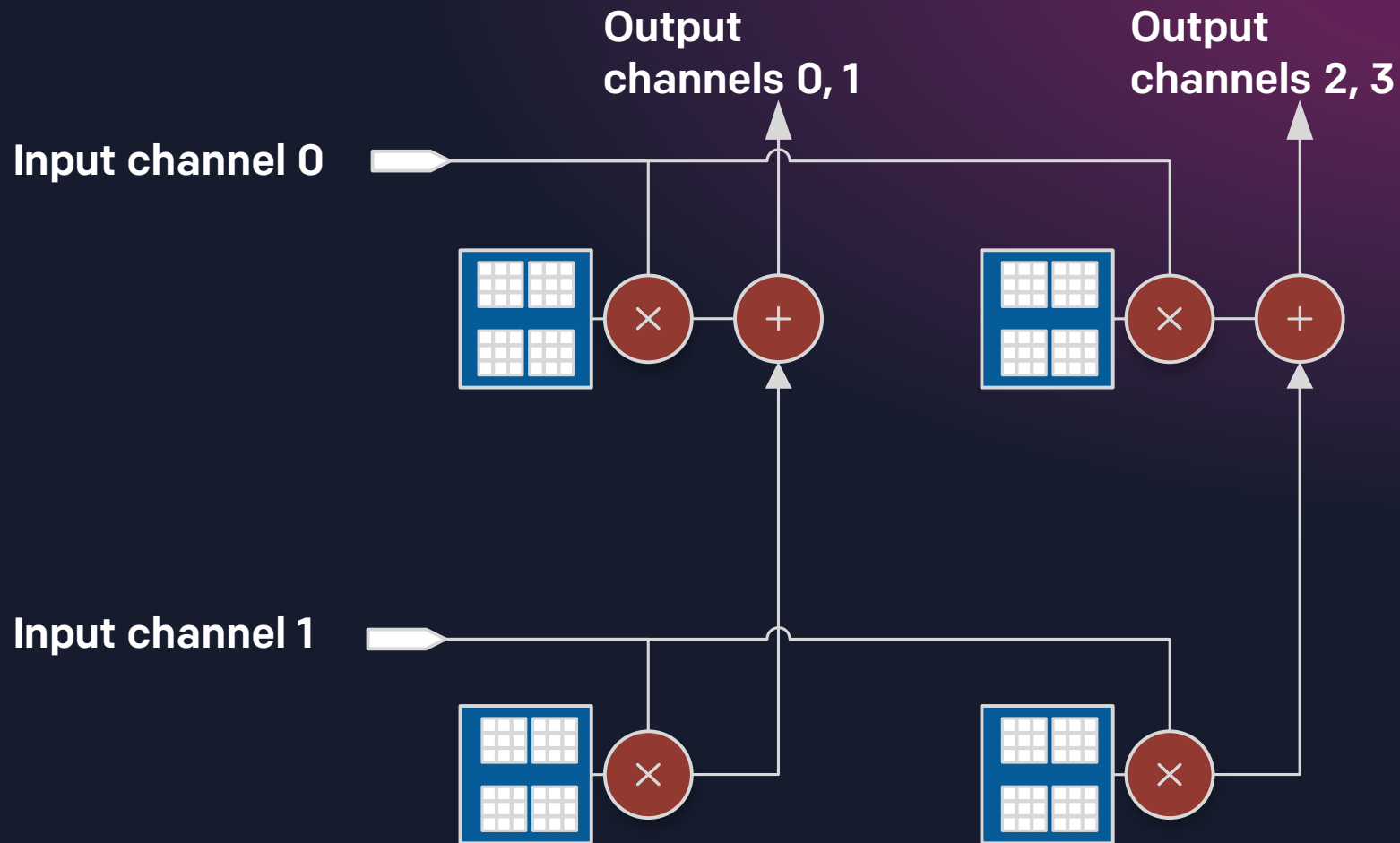
1. Divide weight matrix into blocks.
2. Copy weight block matrices into MxV array.
3. Divide input matrix into column vectors.
4. Run matrix-vector multiply-accumulate.



Maximizing Array Compute Efficiency

Assign input channels to MxV input lanes.

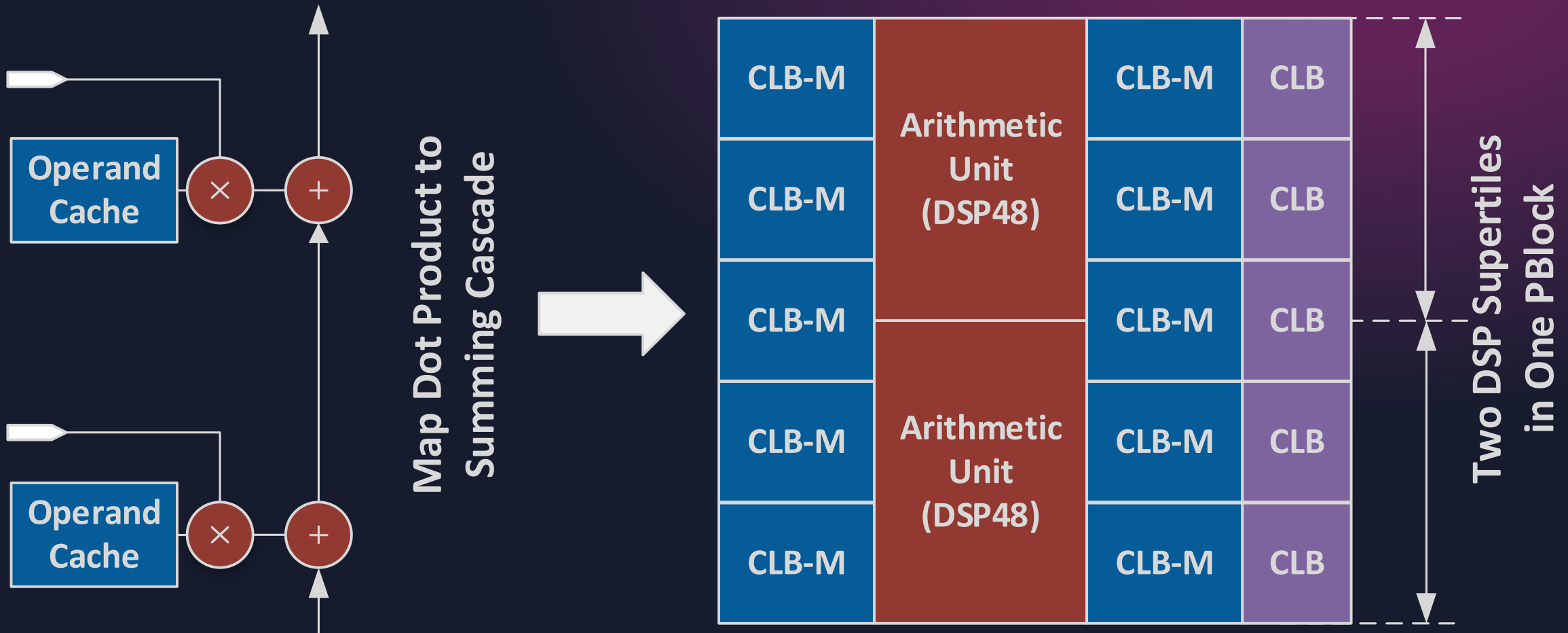
Assign output channels to MxV output lanes.



Maximize Clock Rate: Use DSP Supertiles [FPL '17]

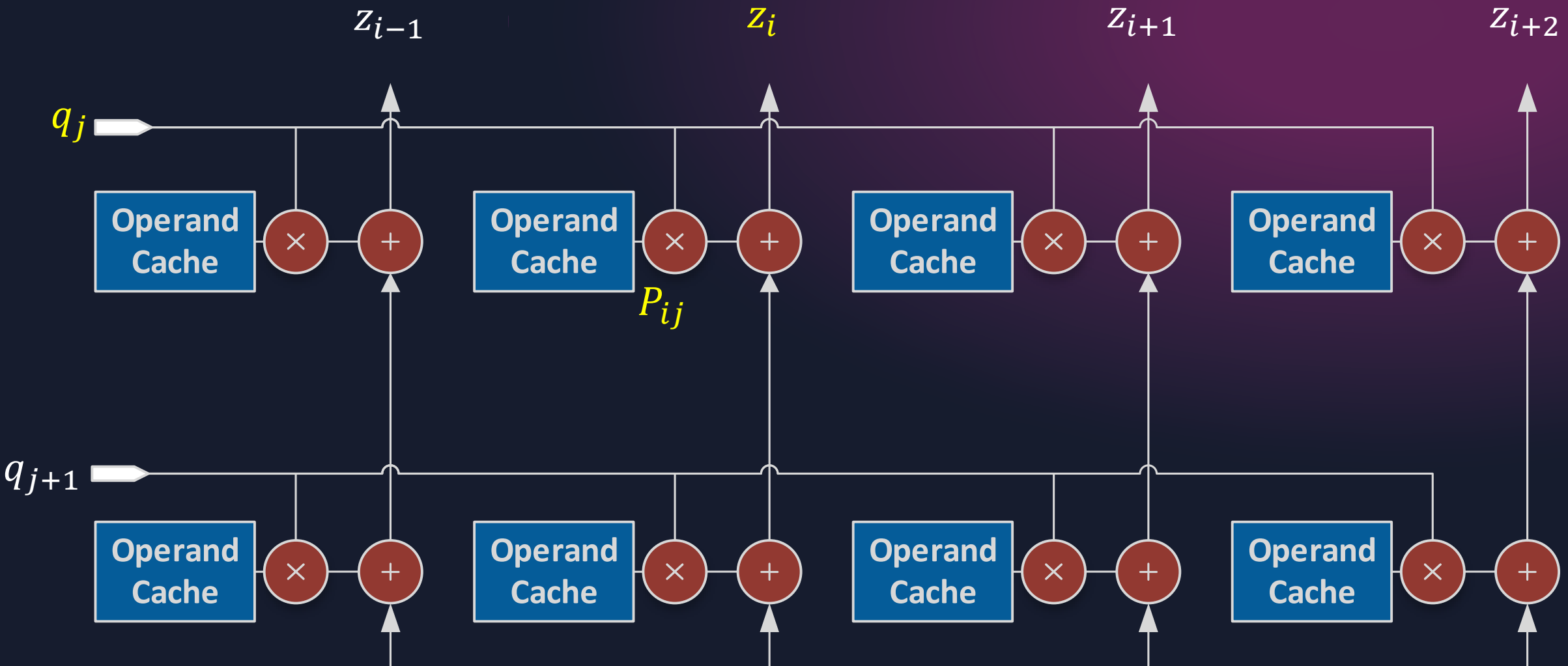
DSP supertile column: 4 physically adjacent columns good for timing closure

Use DSP summing cascades instead of summing trees



Maximize Clock Rate: Use DSP Supertiles

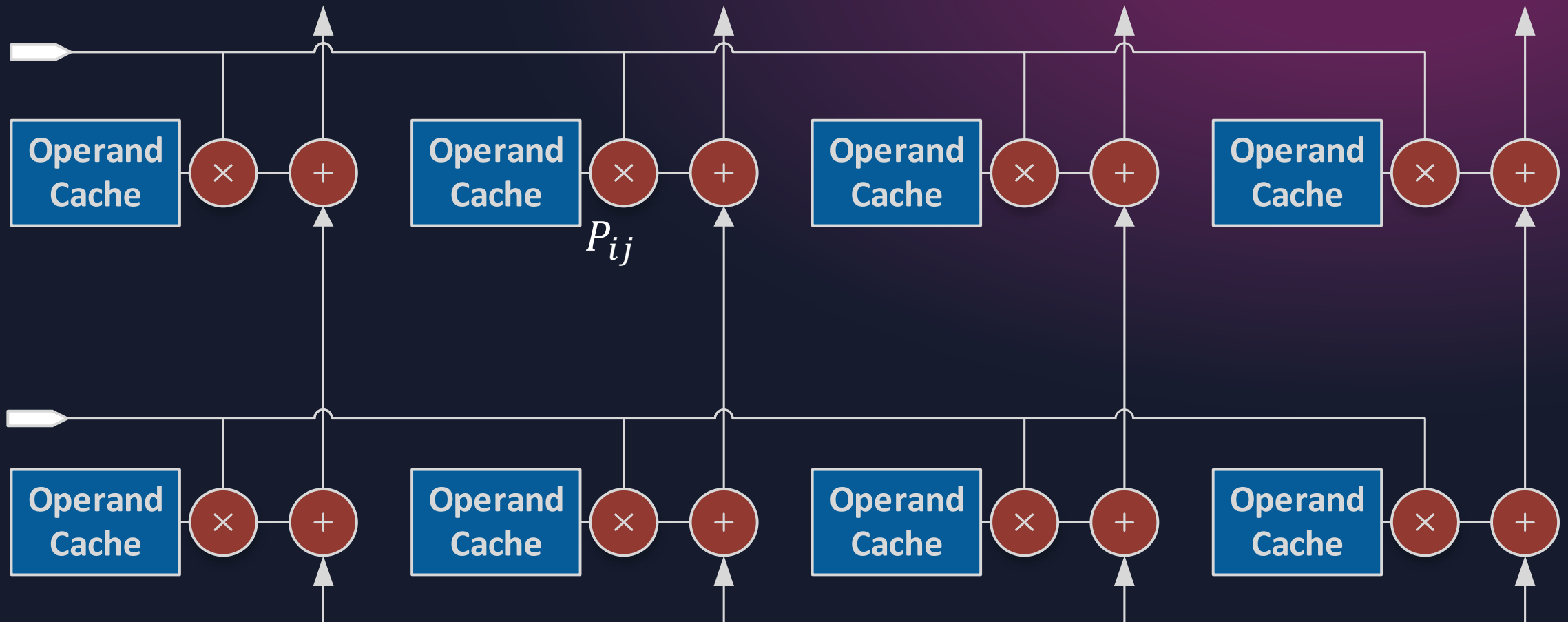
DSP Supertile Array = Matrix-Vector Multiplier: $z \leftarrow Pq + r$



Maximize Clock Rate: Use DSP Supertiles

Double-buffered operand cache in distributed RAM stores filter weights

Distributed RAM can run as fast as DSP tiles



Example: Time Interleaving Four Output Channels

3 × 3 Conv2D Stride 1

Input feature reused
for 4 cycles.

Input channel
address (green box)
moves $\frac{1}{4}$ as fast as
weight address (red
box).



Mapping GoogLeNet to a Processor Chain



GoogLeNet

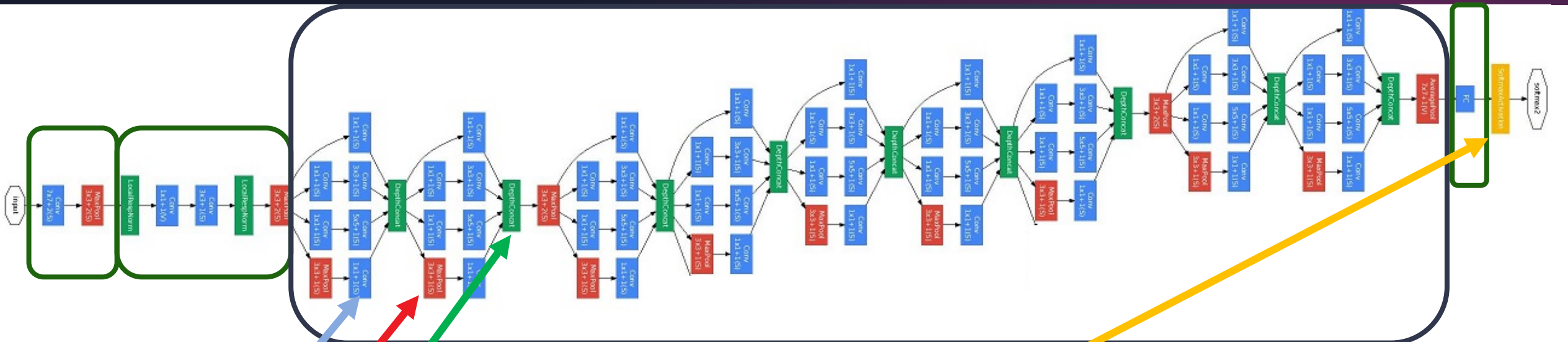
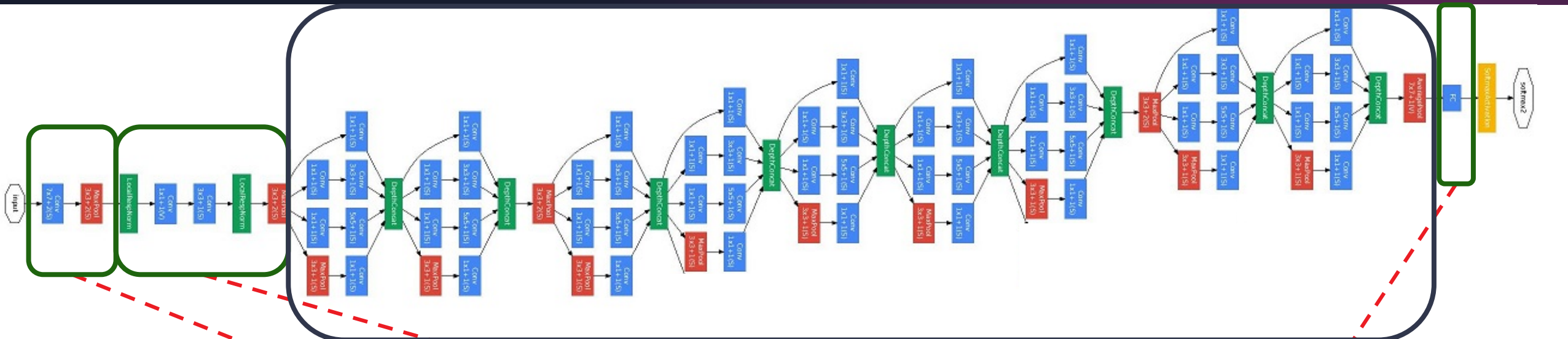


Image in

1000 classes out

- Blue: 2D convolution or fully-connected plus activation (ReLU)
- Red: Max-pooling layer
- Green: Concatenation layer
- Yellow: Softmax

Mapping GoogLeNet to a Four-Processor Chain



+

Input bias

P1

P2

P3

P4

Softmax

One Chain on Each VU9P Super Logic Region (SLR)

There are 3 SLRs on each VU9P.

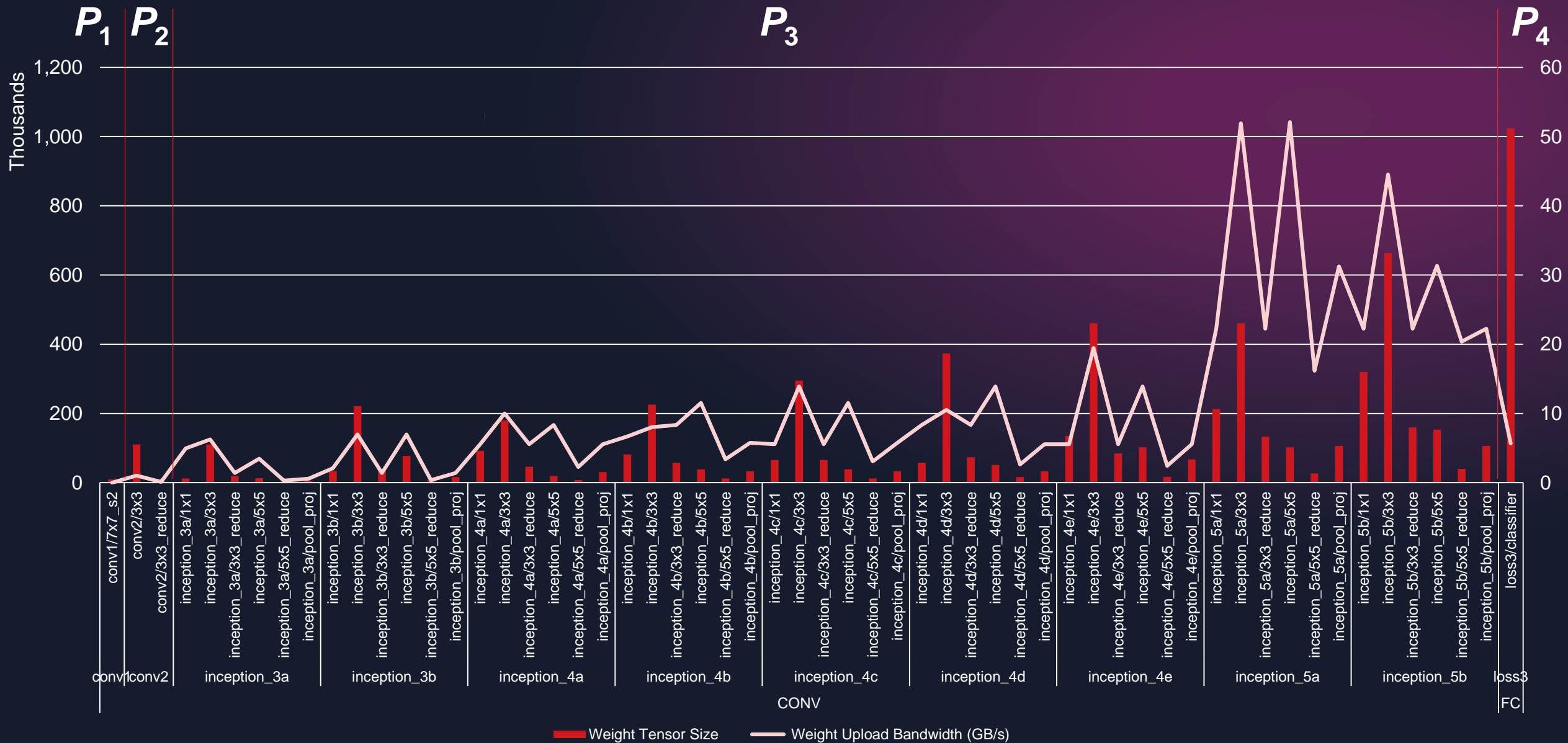
Each SLR contains a 4-processor chain.

Each chain is independent and has its own set of weights.

Each processor in the chain is event-driven (w/o central controller).

Processor	DSP Supertile Array Input Lanes N_1	DSP Supertile Array Output Lanes N_2	Linear Layers in GoogLeNet	Cycles Per Image
P_1	21	8	1	691792
P_2	32	16	2	686432
P_3	96	16	54	708671
P_4	8	1	1	128000

Fast and Big Memory Wanted

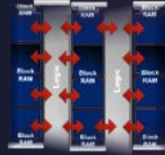


Keep All Tensors On-Die in UltraRAM and BRAM

Bandwidth



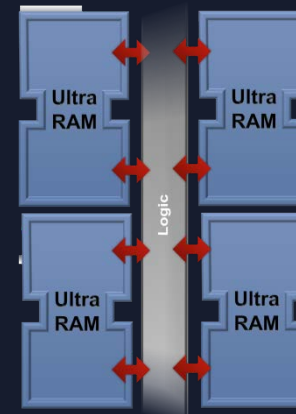
Distributed RAM
(bits to kilobits)



Block RAM
(10s of megabits)

MxV doesn't have to wait for memory

New in UltraScale+™



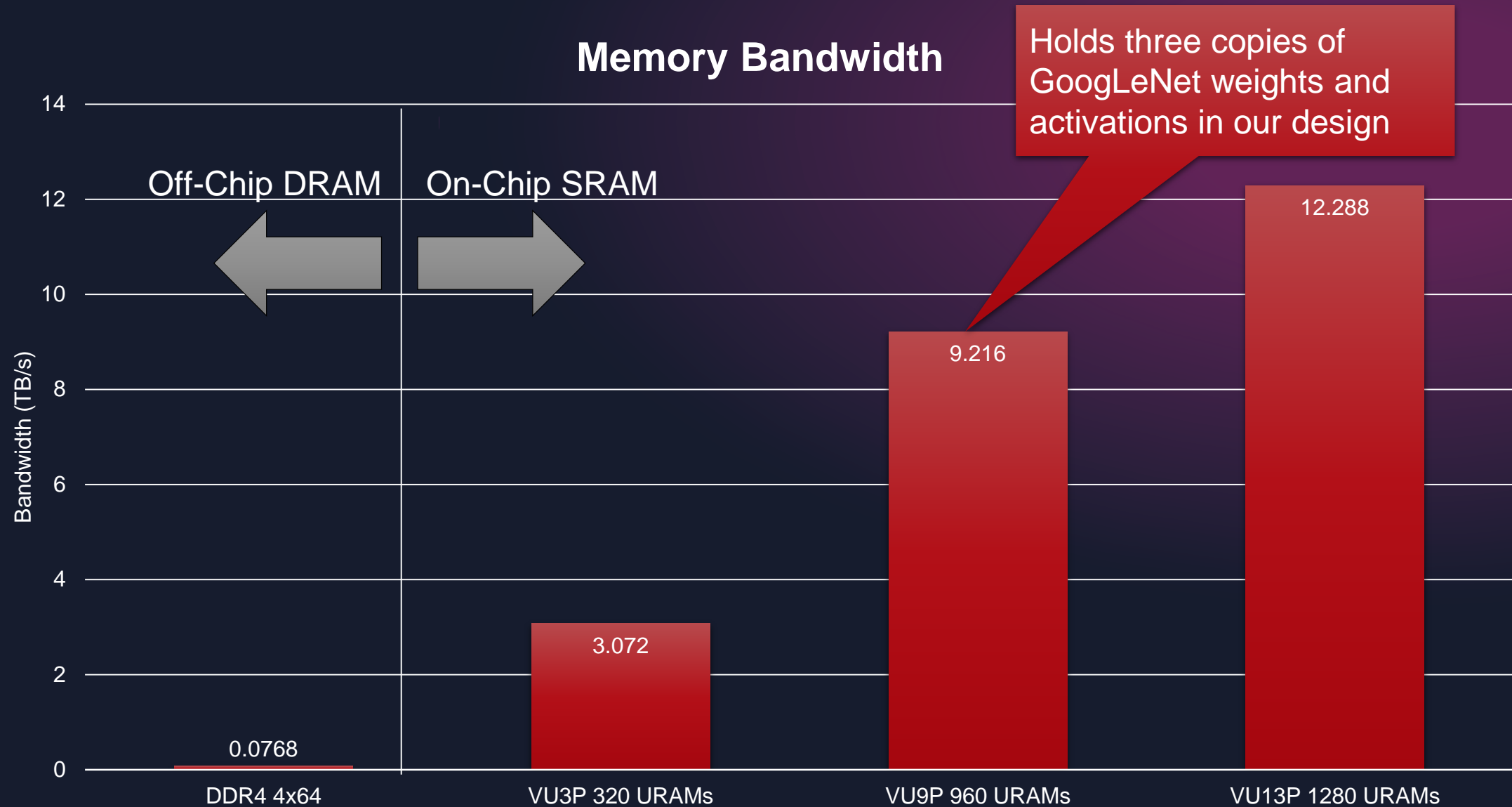
UltraRAM
(100s of megabits)



External DDR DRAM
(10s of gigabits)

Capacity

UltraRAM vs. DDR4 Bandwidth

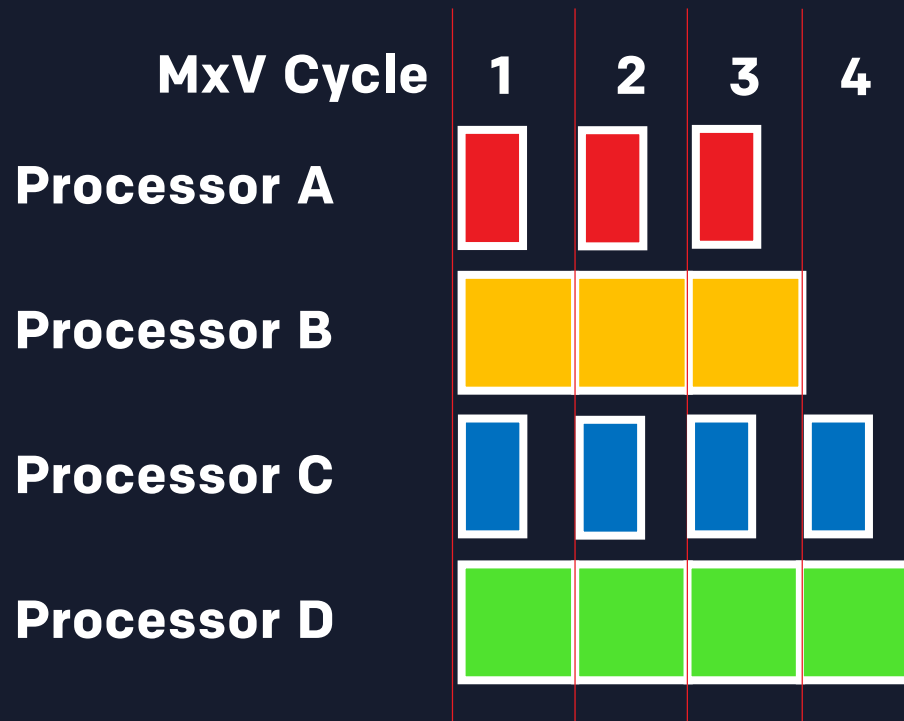


Two Components of Compute Efficiency

Balance factor is a “hurry-up-and-wait” metric.

Distribution efficiency measures multiply-add cycle utilization.

Compute Efficiency = Balance Factor × Distribution Efficiency



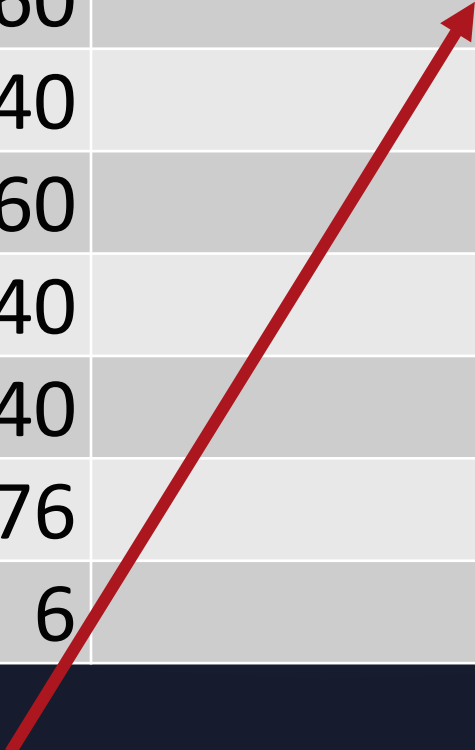
Balance Factor	Distribution Efficiency
Low	Low
Low	High
High	Low
High	High

GoogLeNet Efficiency Metrics (ImageNet 1K Dataset)

Processor	MxV Shape	Cycles Per Image	Balance Factor	Distribution Efficiency	Compute Efficiency
P_1	21×8	691792	97.6%	100.0%	97.6%
P_2	32×16	686432	96.9%	100.0%	96.9%
P_3	96×16	708671	100.0%	95.2%	95.2%
P_4	8×1	128000	18.1%	100.0%	18.1%
Pipeline		708671	98.8%	96.7%	95.5%

Resource Utilization

Resource	Used	Available	Utilization (%)
URAM	960	960	100.0
DSP48E2	3817	6840	55.8
RAMB36E2	791.5	2160	40.4
LUT Memory	129019	591840	21.8
LUT as Logic	288463	1182240	24.4
Serdes	4	76	5.3
PCIE40E4	1	6	16.7



Can close timing despite 100% URAM utilization due to simple SRAM-DSP interconnections

Summary

- > **Make compute fast: Use DSP supertile arrays**
- > **Keep compute busy**
 - >> **Store tensors in SRAM (UltraRAM and BRAM)**
 - >> **Map convolution channels to MxV ports**
 - >> **Configure MxV shapes to maximize pipeline balance factors**
- > **Simplify implementation**
 - >> **Use straight routes. (Recall the flat view.)**
 - >> **Don't build convoluted memory -compute networks**
 - >> **Use temporal locality to reduce tensor memory read bandwidth**
 - >> **Run the rest of the design at half the DSP clock rate**