



arm

Efficient Winograd or Cook-Toom Convolution Kernel Implementation on Widely Used Mobile CPUs

Partha Maji, Andrew Mundy, Ganesh Dasika,
Jesse Beu, Matthew Mattina, Robert Mullins

Arm Research, University of Cambridge

ML and the Rise of the Edge

VR/AR/MR



Robotics



Drones



Shipping & logistics



IoT



Home, surveillance & analytics



Automotive



Mobile



Contributions of this work

- We discuss what Winograd convolution can offer in terms of performance
- Breakdown the instruction-level implications and memory layout tradeoffs for different flavors of a Winograd kernel in order to realize its full potential
- Demonstrate how general matrix multiply (GEMM) can further optimize Winograd
- Present performance results for Winograd vs conventional im2row + GEMM solution
 - More than a 2x performance boost on real hardware today!

Ultimately enable more efficient ML compute at the edge through Winograd in the Arm Compute Library (ArmCL).

arm

Convolution and
Winograd

What is Winograd and why should I care?

- Convolutional Neural Networks (CNNs)
 - Common type of deep learning model employed in a variety of domains
 - Convolve filter bank (weights) over a field (input activations) to produce a response (output)
 - Push response through an activation function (typically ReLu) and feed to the next layer
- Winograd Convolution
 - Based in the Chinese Remainder Theorem and modulo arithmetic
 - Produces mathematically equivalent results to naïve convolution*
 - Similar to using Fourier: transform into 'Winograd domain', do simpler math, transform result back

*Assuming infinite precision

What is Winograd and why should I care?

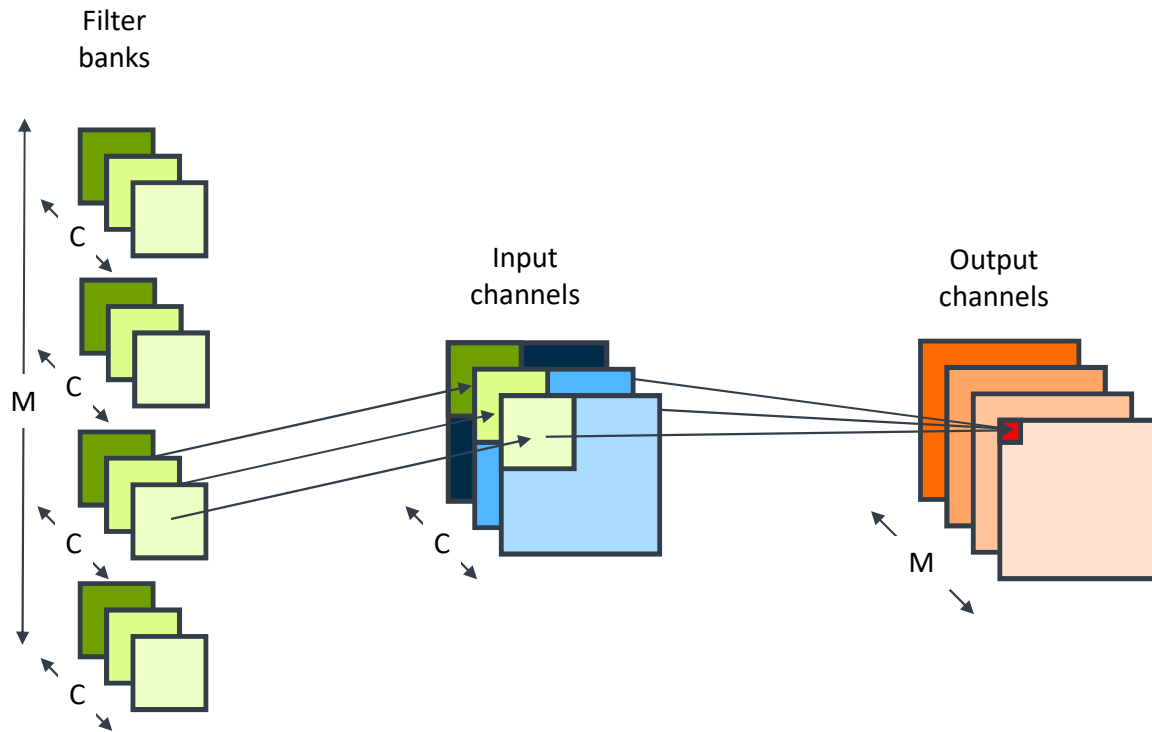
- Convolutional Neural Networks (CNNs)
 - Common type of deep learning model employed in a variety of domains
 - Convolve filter bank (weights) over a field (input activations) to produce a response (output)
 - Push response through an activation function (typically ReLu) and feed to the next layer
- Winograd Convolution
 - Based in the Chinese Remainder Theorem and modulo arithmetic
 - Produces mathematically equivalent results to naïve convolution*
 - Similar to using Fourier: transform into 'Winograd domain', do simpler math, transform result back

Objective: To (quickly) explain for a CPU context:

$$f = Z^T [(WwW^T) \odot (X^T x X)] Z$$

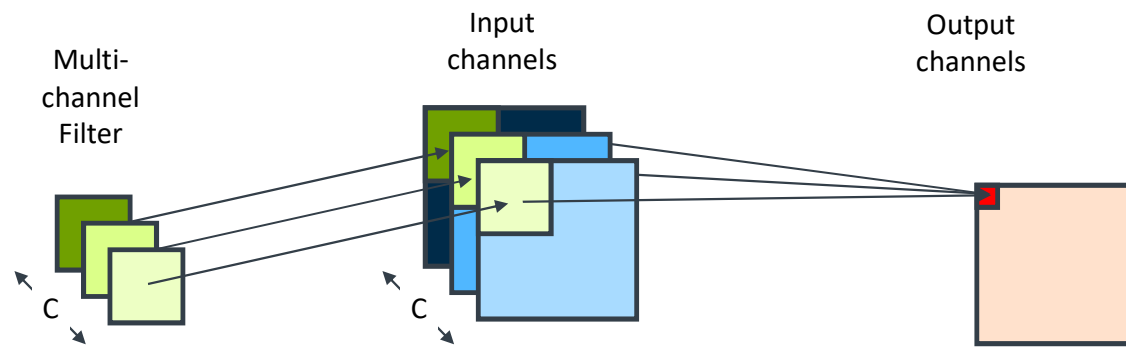
*Assuming infinite precision

Winograd Convolution

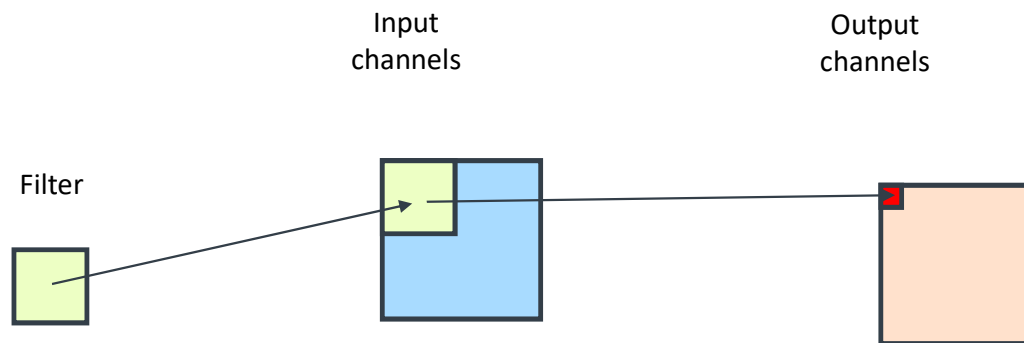


Standard CNN Configuration

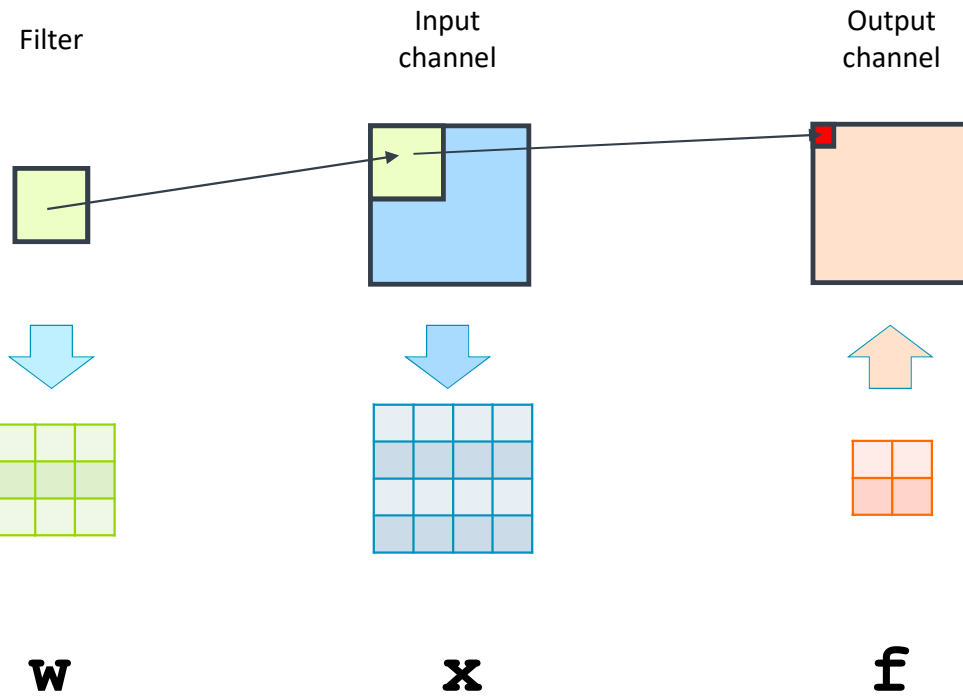
Winograd Convolution



Winograd Convolution



Winograd Convolution



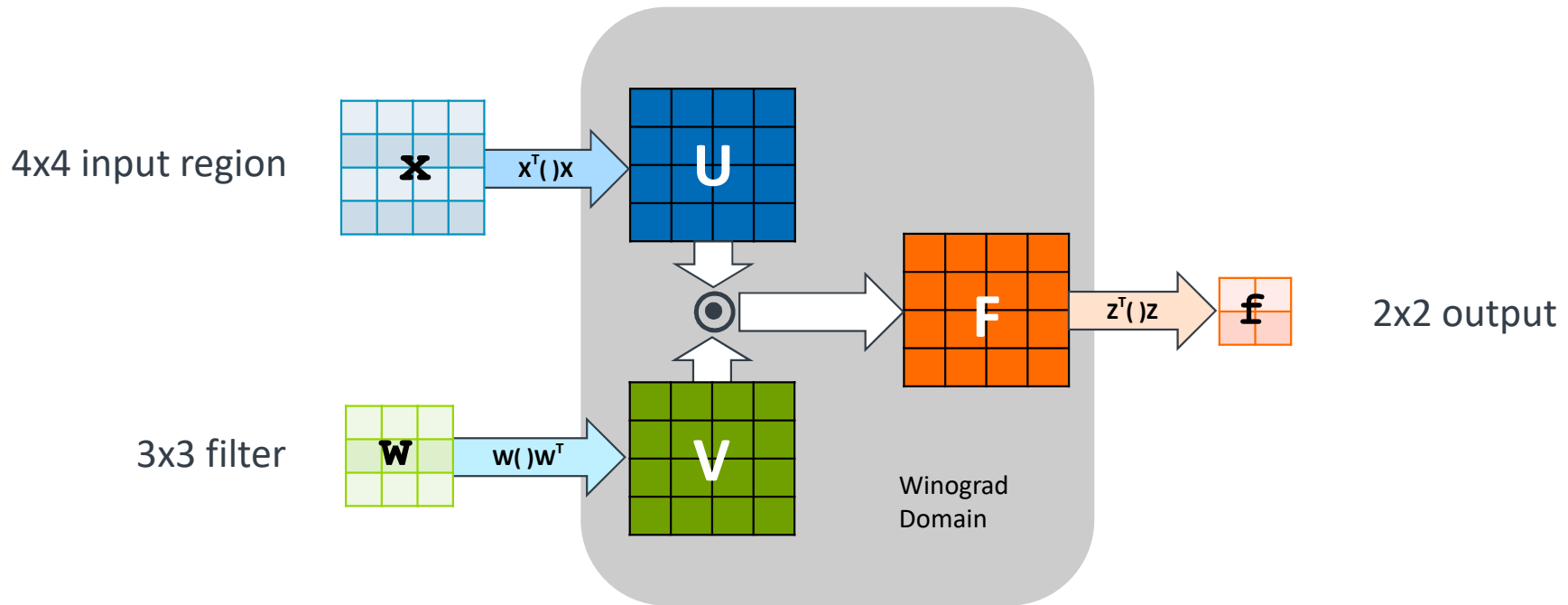
Assume:

3x3 filter

4x4 chunk of input activations
(aka, a 'region')

2x2 output

Winograd Convolution



$$\mathbf{f} = \mathbf{Z}^T \left[(\mathbf{W}\mathbf{w}\mathbf{W}^T) \odot (\mathbf{X}^T \mathbf{x} \mathbf{X}) \right] \mathbf{Z}$$

Input Region Transform

$$(2 \times 2) = (2 \times 4) [(4 \times 3)(3 \times 3)(3 \times 4) \odot (4 \times 4)(4 \times 4)(4 \times 4)] (2 \times 4)$$

$$\left[\begin{array}{c} \left[\begin{array}{|c|c|c|c|} \hline 1 & 0 & -1 & 0 \\ \hline 0 & 1 & 1 & 0 \\ \hline 0 & -1 & 1 & 0 \\ \hline 0 & 1 & 0 & -1 \\ \hline \end{array} \right] \cdot \left[\begin{array}{|c|c|c|c|} \hline \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square \\ \hline \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square \\ \hline \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square \\ \hline \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square & \color{lightblue} \square \\ \hline \end{array} \right] \cdot \left[\begin{array}{|c|c|c|c|} \hline 1 & 0 & 0 & 0 \\ \hline 0 & 1 & -1 & 1 \\ \hline -1 & 1 & 1 & 0 \\ \hline 0 & 0 & 0 & -1 \\ \hline \end{array} \right] \\ \mathbf{X}^T \qquad \qquad \mathbf{x} \qquad \qquad \mathbf{X} \\ \end{array} \right] = \mathbf{U}$$

$$f = \mathbf{Z}^T [(WwW^T) \odot (\mathbf{X}^T \mathbf{x} \mathbf{X})] \mathbf{Z}$$

Filter Transform

$$(2 \times 2) = (2 \times 4) \left[(4 \times 3)(3 \times 3)(3 \times 4) \odot (4 \times 4)(4 \times 4)(4 \times 4) \right] (2 \times 4)$$

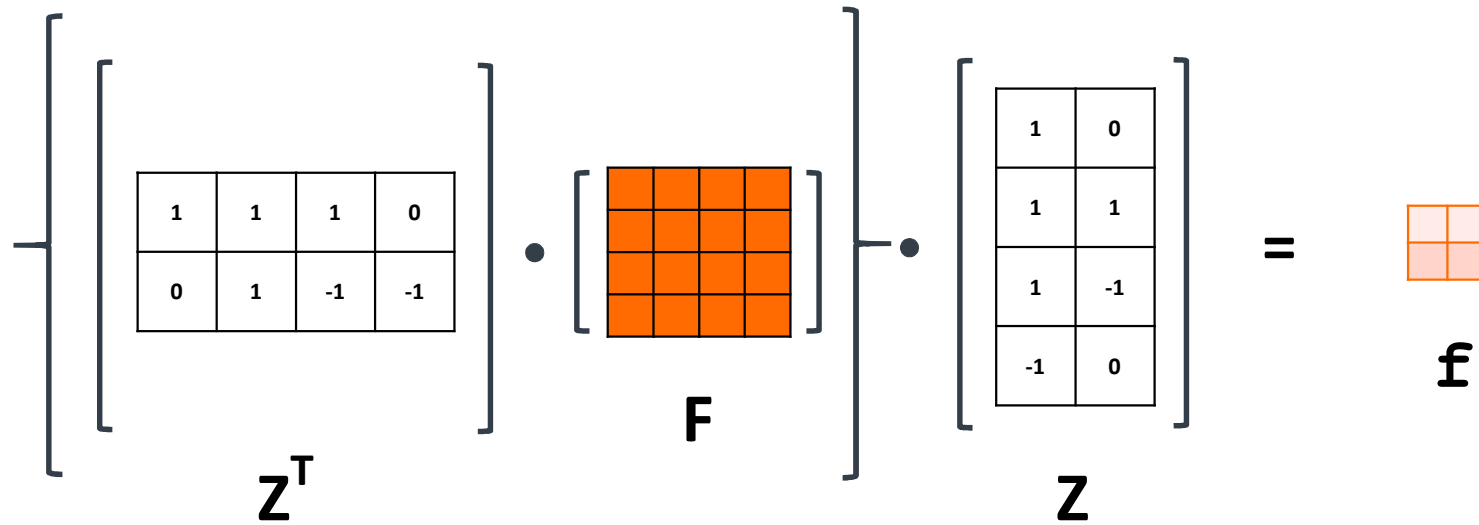
$$\left[\begin{matrix} \begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 1/2 & 1/2 \\ 1/2 & -1/2 & 1/2 \\ 0 & 0 & 1 \end{bmatrix} \\ \mathbf{W} \end{matrix} \right] \cdot \begin{matrix} \begin{bmatrix} & & \\ & & \\ & & \\ & & \end{bmatrix} \\ \mathbf{w} \end{matrix} \right] \cdot \begin{matrix} \begin{bmatrix} 1 & 1/2 & 1/2 & 0 \\ 0 & 1/2 & -1/2 & 0 \\ 0 & 1/2 & 1/2 & 1 \end{bmatrix} \\ \mathbf{W}^T \end{matrix} \right] = \begin{matrix} \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \\ \mathbf{V} \end{matrix}$$

$$f = Z^T \left[(\mathbf{W} \mathbf{w} \mathbf{W}^T) \odot (\mathbf{X}^T \mathbf{x} \mathbf{X}) \right] \mathbf{Z}$$

Output Channel Transform

$$(2 \times 2) = (2 \times 4) \left[(4 \times 3)(3 \times 3)(3 \times 4) \odot (4 \times 4)(4 \times 4)(4 \times 4) \right] (4 \times 2)$$

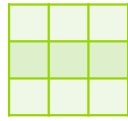
(this reduces down to a 4x4)



$$f = Z^T \left[(WwW^T) \odot (X^T x X) \right] Z$$

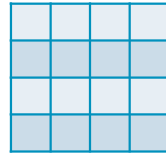
Elementwise Multiplication

Spatial
Domain



W

\otimes



X

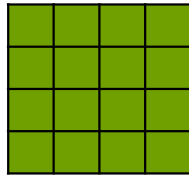
=



F

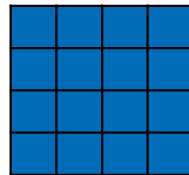
36 MACs

Winograd
Domain



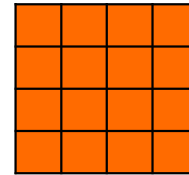
U

\odot



V

=

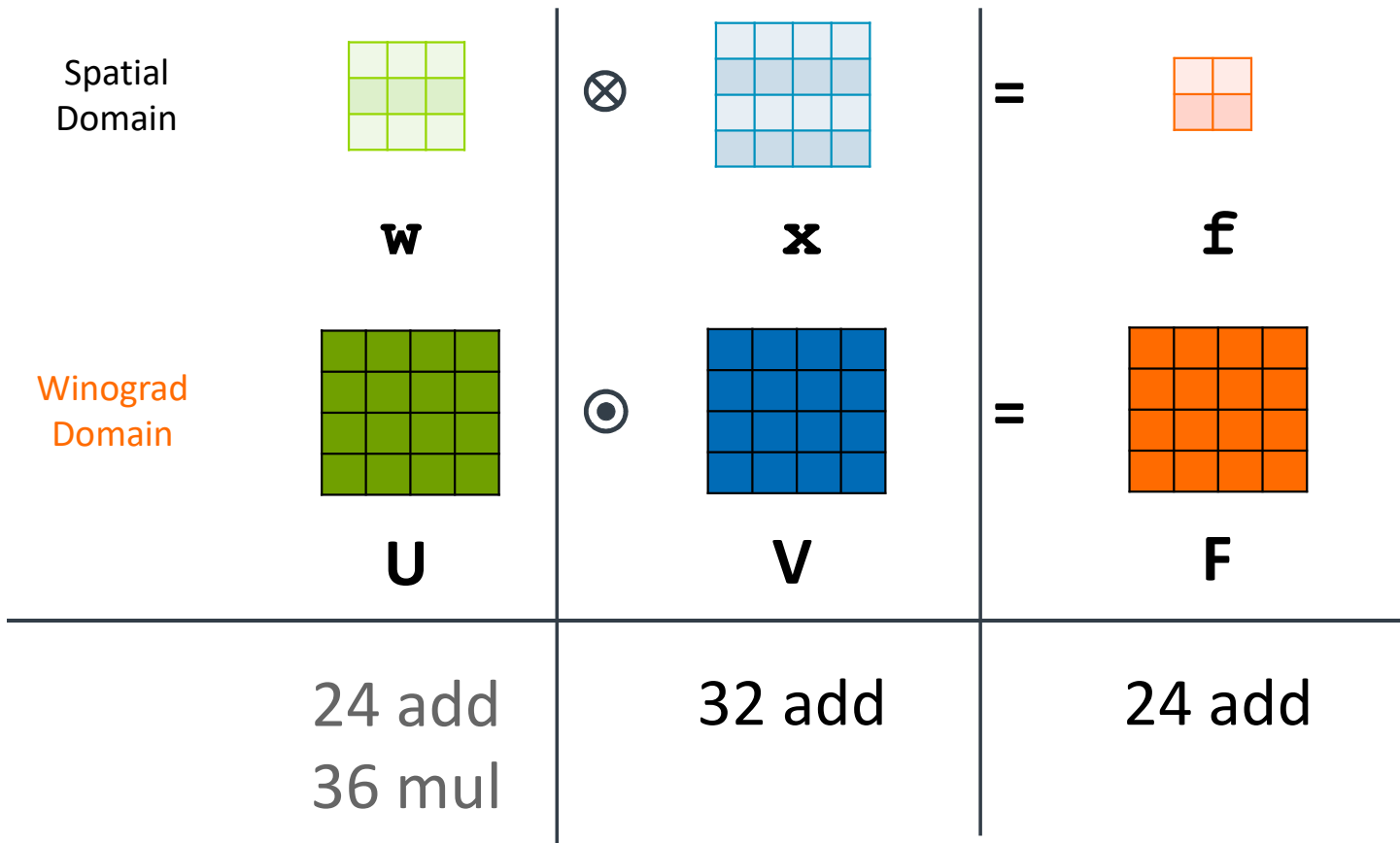


F

16 mul

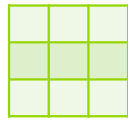
$36 / 16 = 2.25x$ reduction in ops

Transform Cost



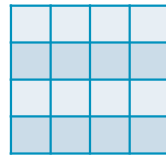
Transform Cost

Spatial Domain



W

\otimes



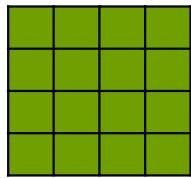
X

=



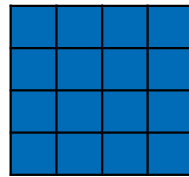
f

Winograd Domain



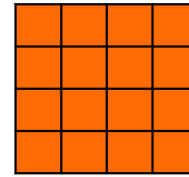
U

\odot



V

=



F

Offline Cost

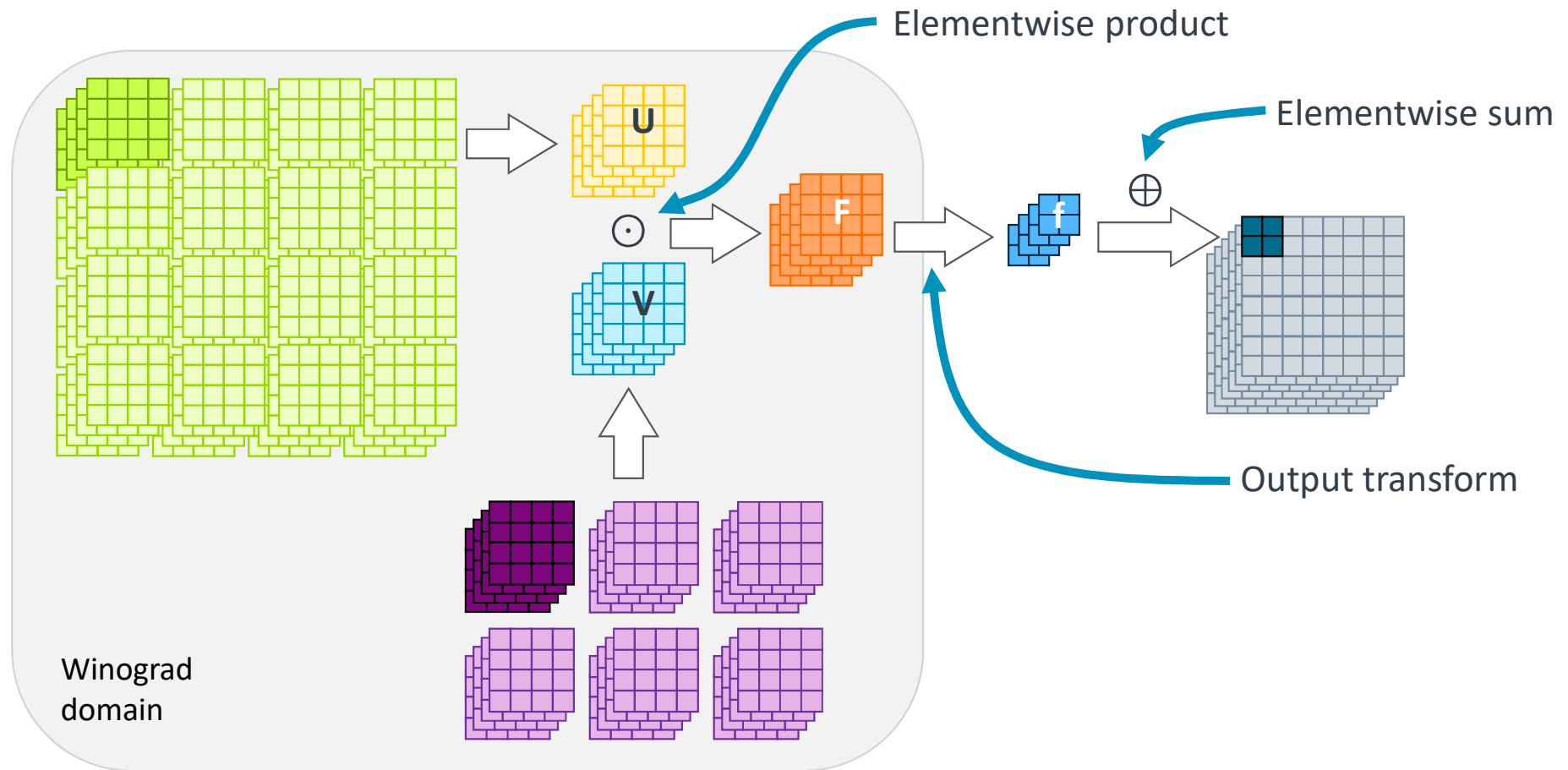
24 add
36 mul

32 add
One-time cost,
high reuse

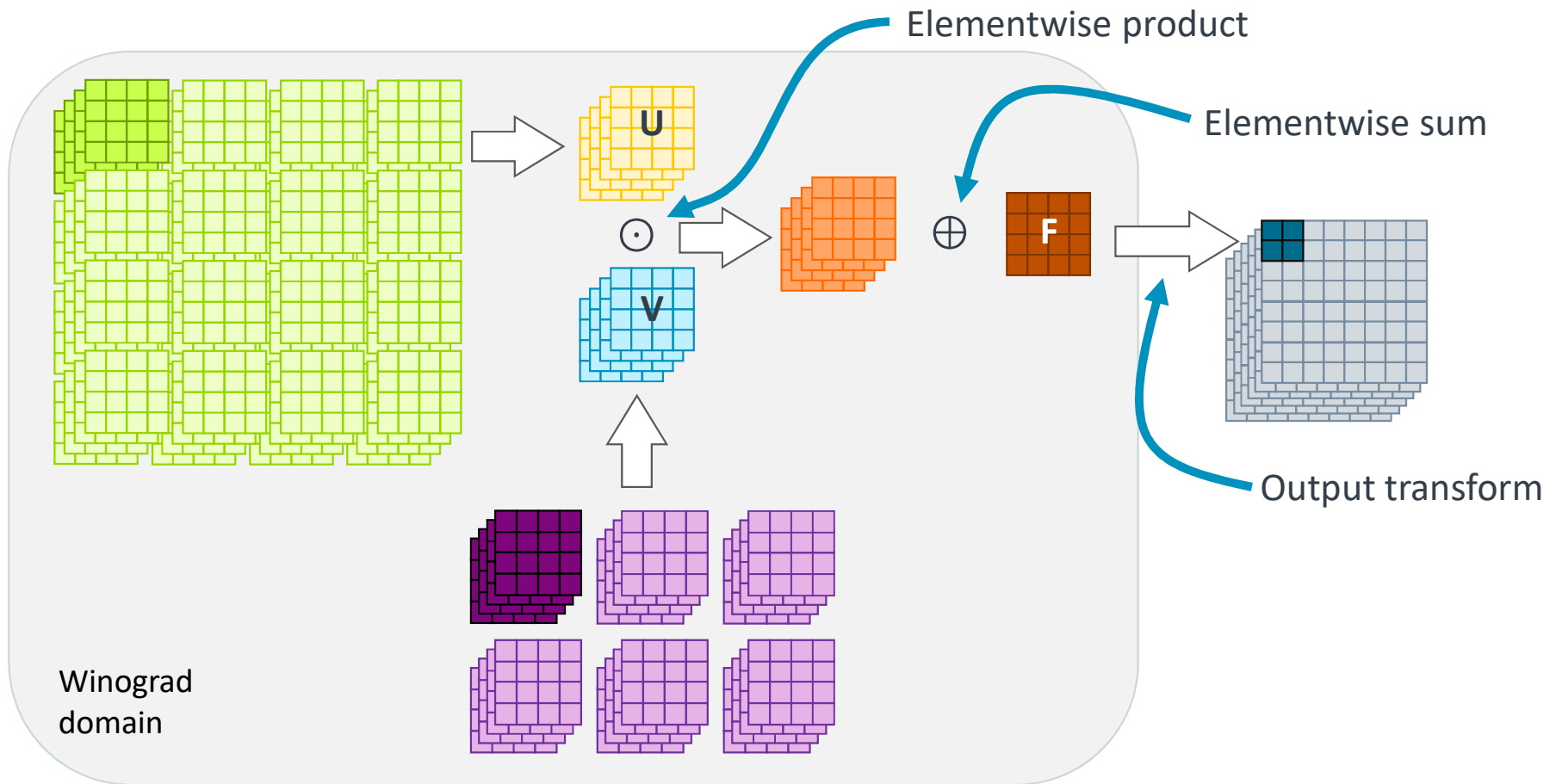
24 add

What about this??
Go Multi-channel!

Winograd for multichannel convolution



Winograd for multichannel convolution – rearranged



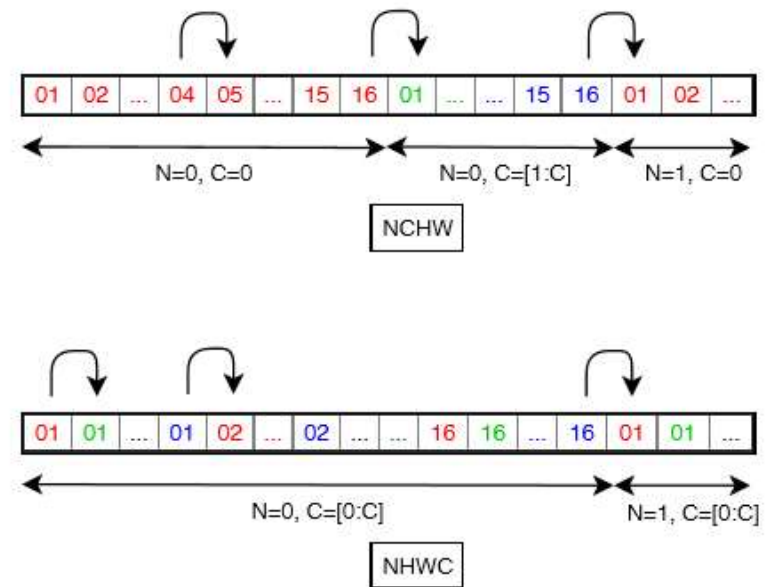
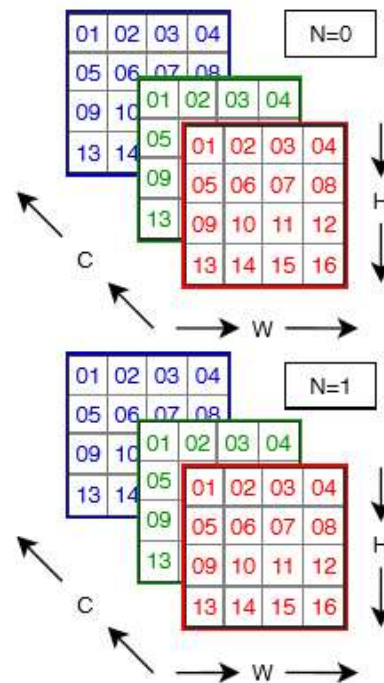
arm

Multi-Channel Filters, Memory Layout, Vectorization, and GEMM

NCHW vs NHWC, data layout

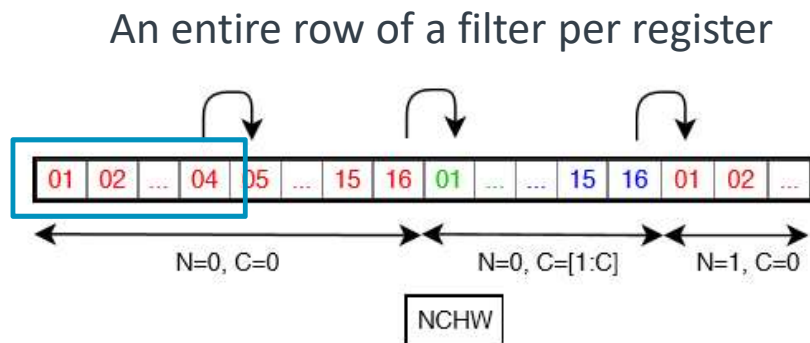
Tensor Ordering

- N = batch
- C = channel
- H = height
- W = width

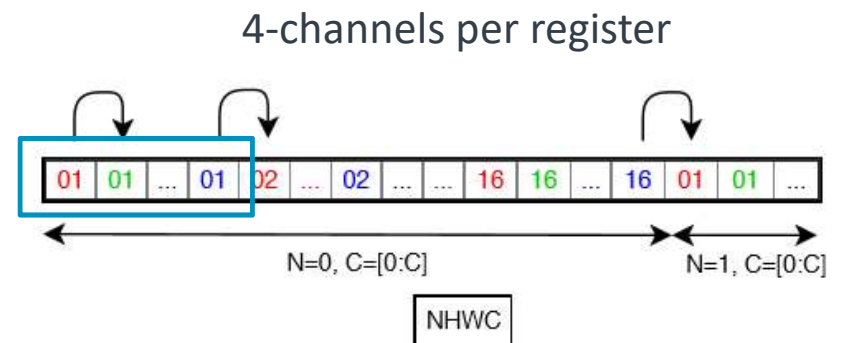


NCHW vs NHWC, data layout

- Layout ultimately dictates how contiguous vector-load operations will populate registers
 - Under NCHW, registers will be filled entirely from a single channel
 - Under NHWC, registers will hold multiple channels for a single coordinate
- In the Arm-V8 architecture (with 128-bit SIMD registers), this means either:



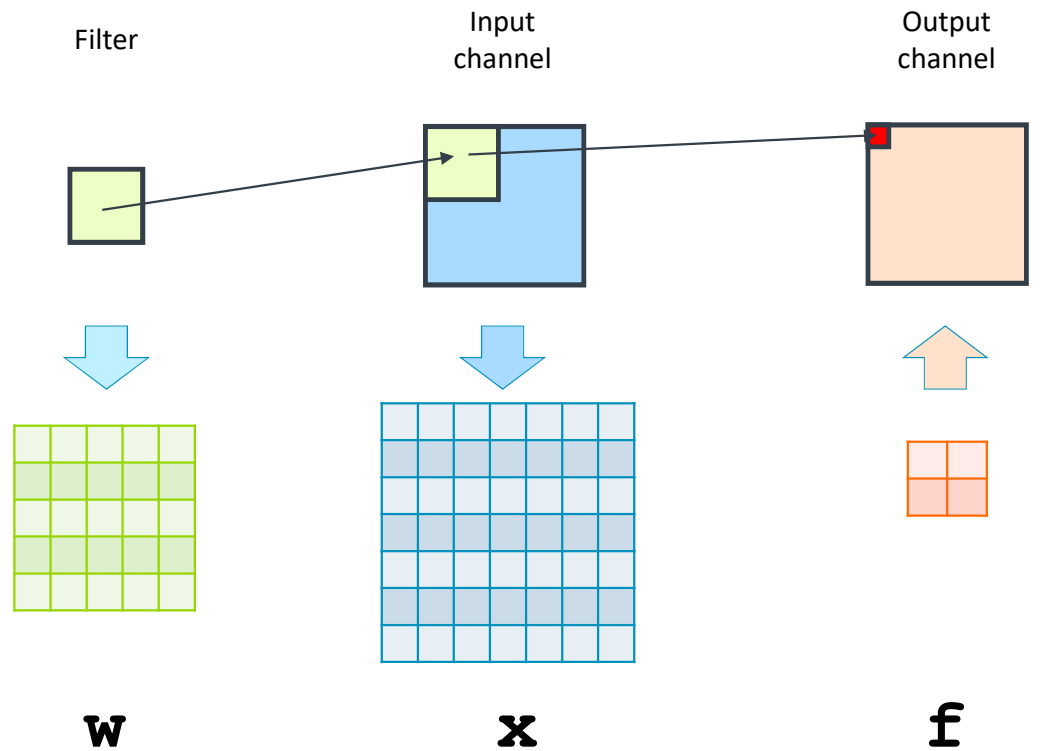
or



Advantages to NHWC layout for CPUs

- Reasonably optimized transforms exist for both NCHW and NHWC at $F(2 \times 2, 3 \times 3, 4 \times 4)$
- Convolution filters and Winograd are not restricted to $F(2 \times 2, 3 \times 3, 4 \times 4)$
 - Larger regions yields can drive higher performance – e.g., $F(3 \times 3, 3 \times 3, 5 \times 5)$
 - 5×5 and 7×7 filters found in inception networks – e.g., $F(2 \times 2, 5 \times 5, 6 \times 6)$
 - Dimension-to-register capacity mismatch results in wasted register utilization and/or alignment complexity under NCHW
 - NHWC only experiences increased register pressure

F(2x2, 5x5, 6x6) Example



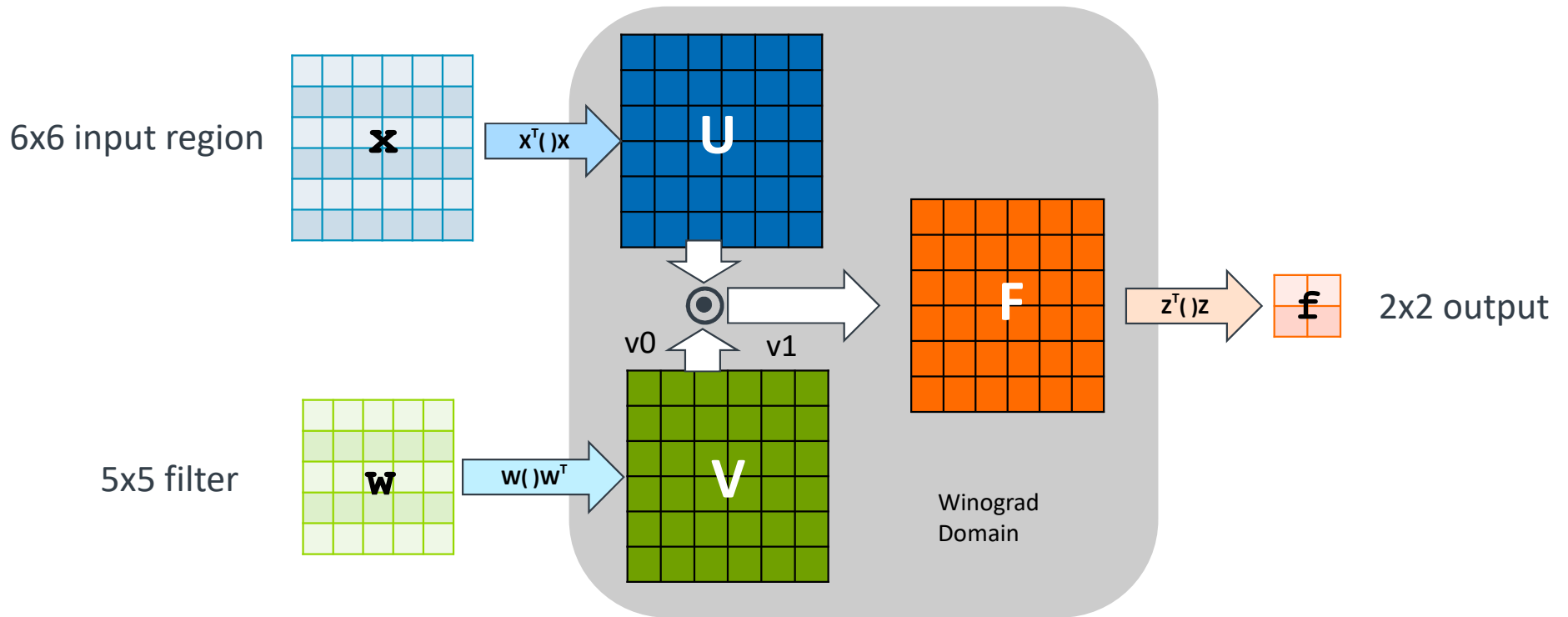
Assume:

5x5 filter

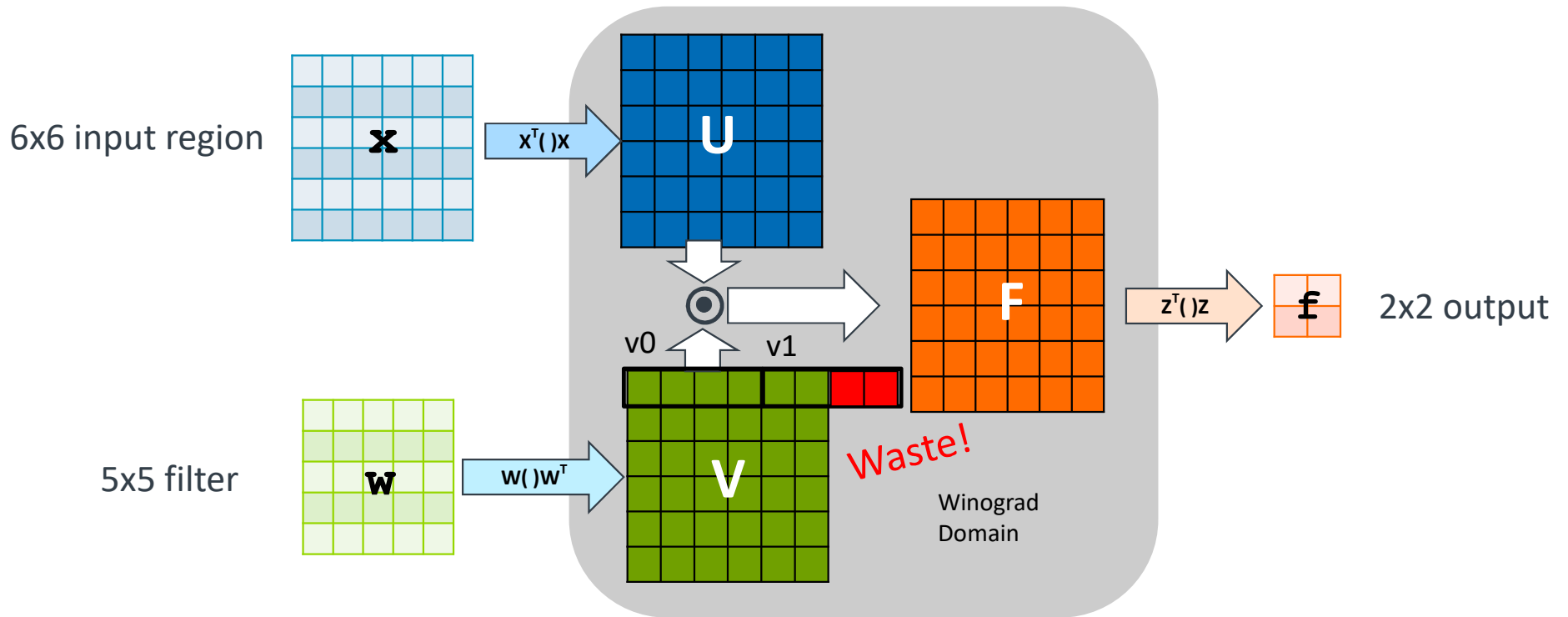
6x6 chunk of input activations

2x2 output

F(2x2, 5x5, 6x6) Example



F(2x2, 5x5, 6x6) Example



Advantages to NHWC layout for CPUs

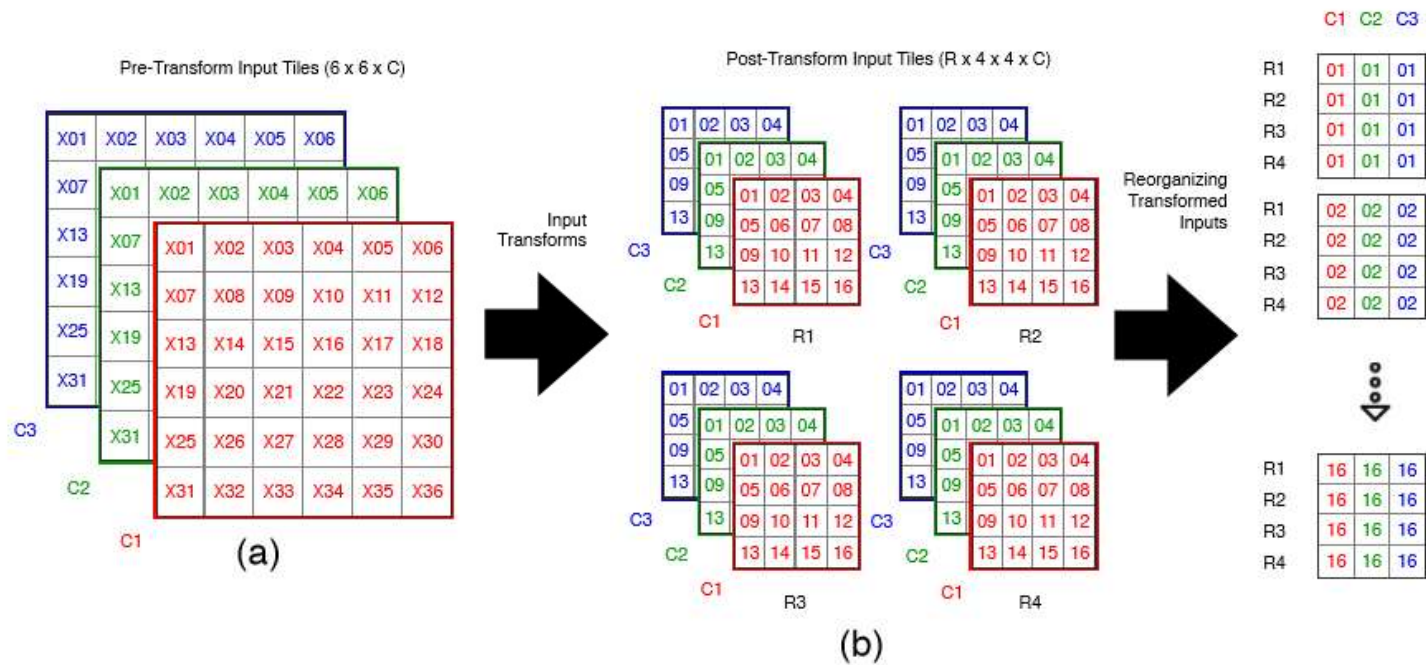
- Reasonably optimized transforms exist for both NCHW and NHWC at $F(2 \times 2, 3 \times 3, 4 \times 4)$
- Convolution filters and Winograd are not restricted to $F(2 \times 2, 3 \times 3, 4 \times 4)$
 - Larger regions yields can drive higher performance – e.g., $F(4 \times 4, 3 \times 3, 6 \times 6)$
 - 5×5 and 7×7 filters found in inception networks – e.g., $F(2 \times 2, 5 \times 5, 7 \times 7)$
 - Dimension-to-register capacity mismatch results in wasted register utilization and alignment complexity under NCHW
 - NHWC only experiences increased register pressure
- Wider registers or lower precision also adds challenges for NCHW
 - 256-bit or FP16 means 8 values per register, or 2 rows per register under NCHW
 - Loss of 1:1 register-row mapping complicates assembly sequence for efficient NCHW transpose
 - NHWC simply doubles the # of channels stored per register

Vectorization over channels is more portable and performant!

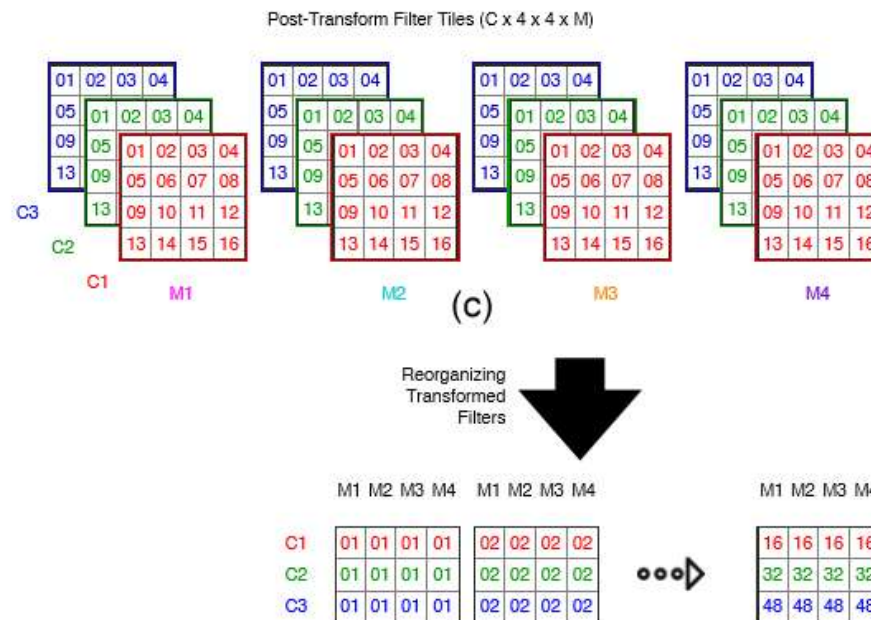
Use of GEMM to further optimize

- General Matrix-Matrix Multiply is a common, highly optimized operation for most architectures, including Arm
- Inspection of the full Winograd convolution algorithm (Listing 1 in paper) shows:
 - The fundamental operation is a multiply-accumulate
 - There are 2 axis of data re-use:
 - weight tile reuse over all input regions and
 - input region reuse over all output channels
 - Opportunity to leverage GEMM to do the computation in a highly parallel manner

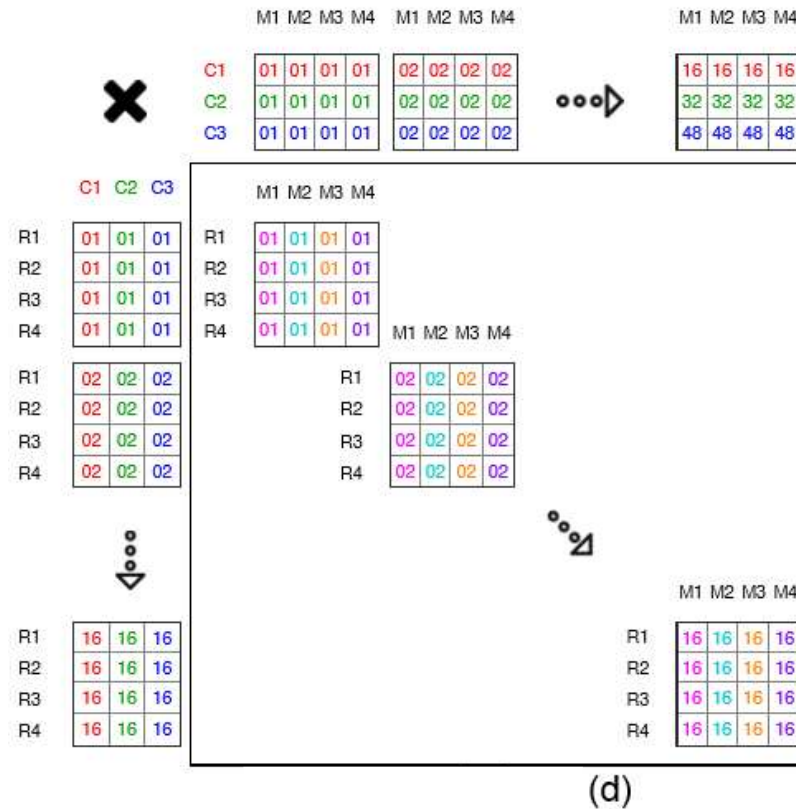
Winograd execution using Matrix of GEMMs



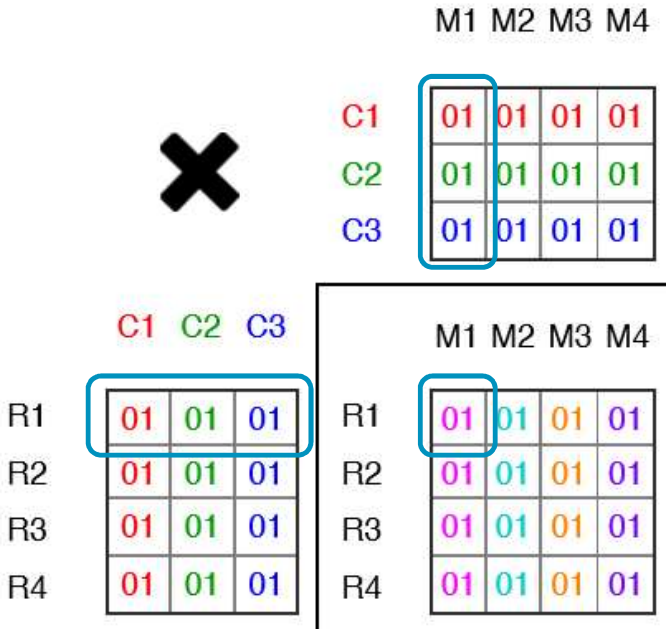
Winograd execution using Matrix of GEMMs



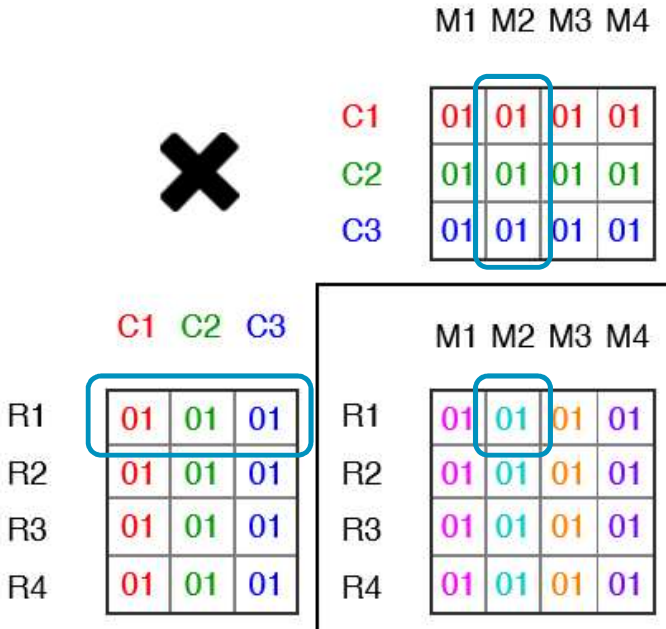
Winograd execution using Matrix of GEMMs



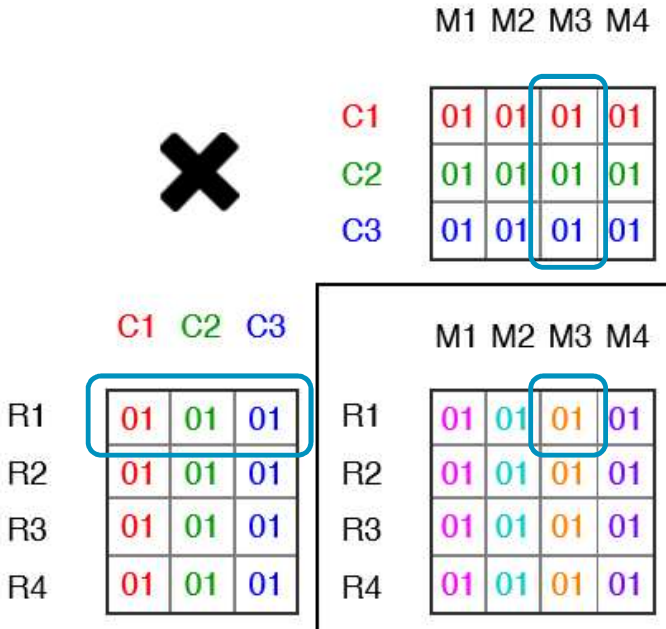
Zoom on individual GEMM



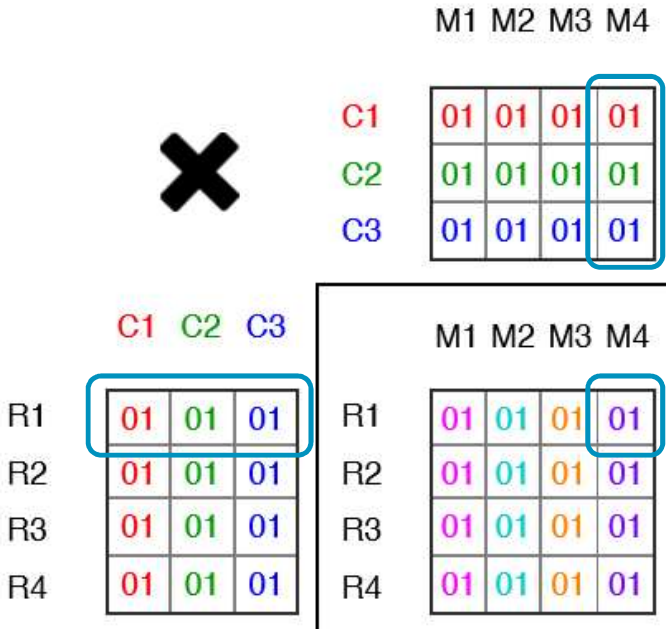
Zoom on individual GEMM



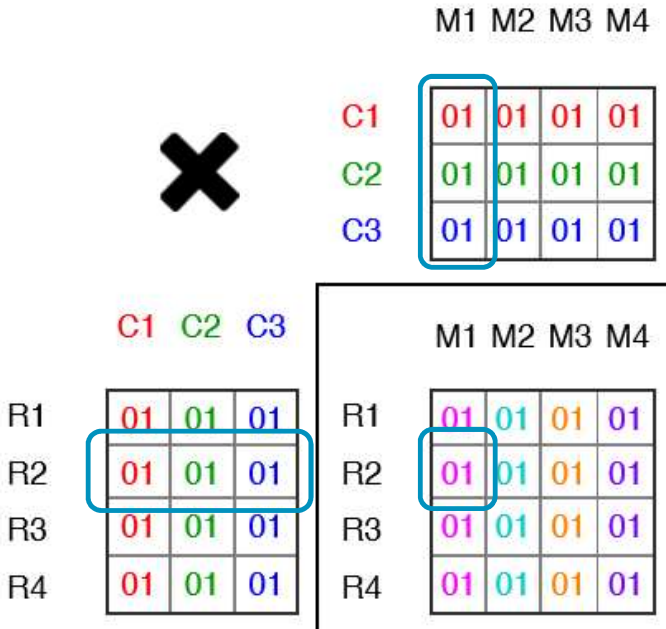
Zoom on individual GEMM



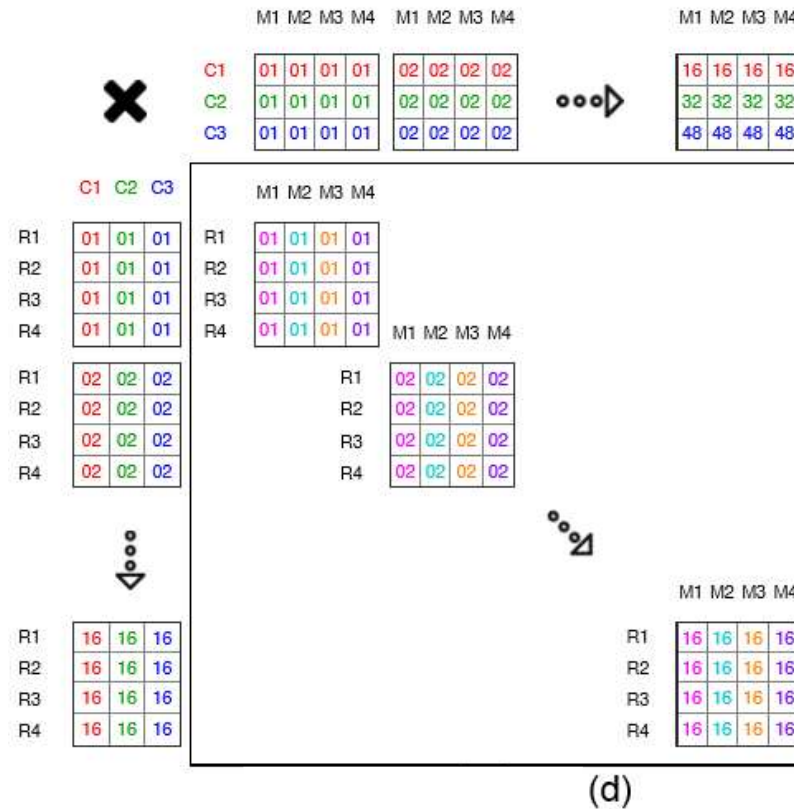
Zoom on individual GEMM



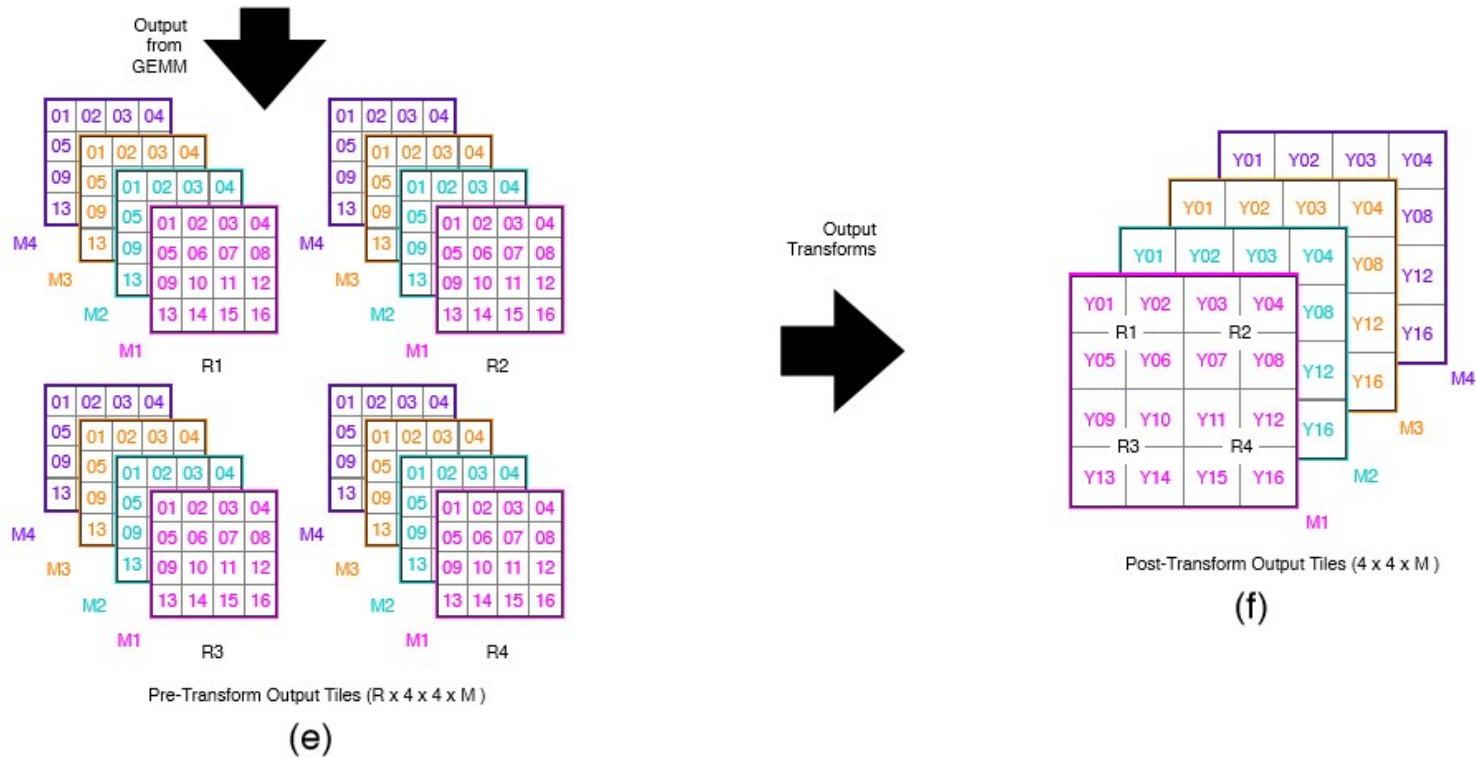
Zoom on individual GEMM



Winograd execution using Matrix of GEMMs



Winograd execution using Matrix of GEMMs



arm

Results

Experimental Setup

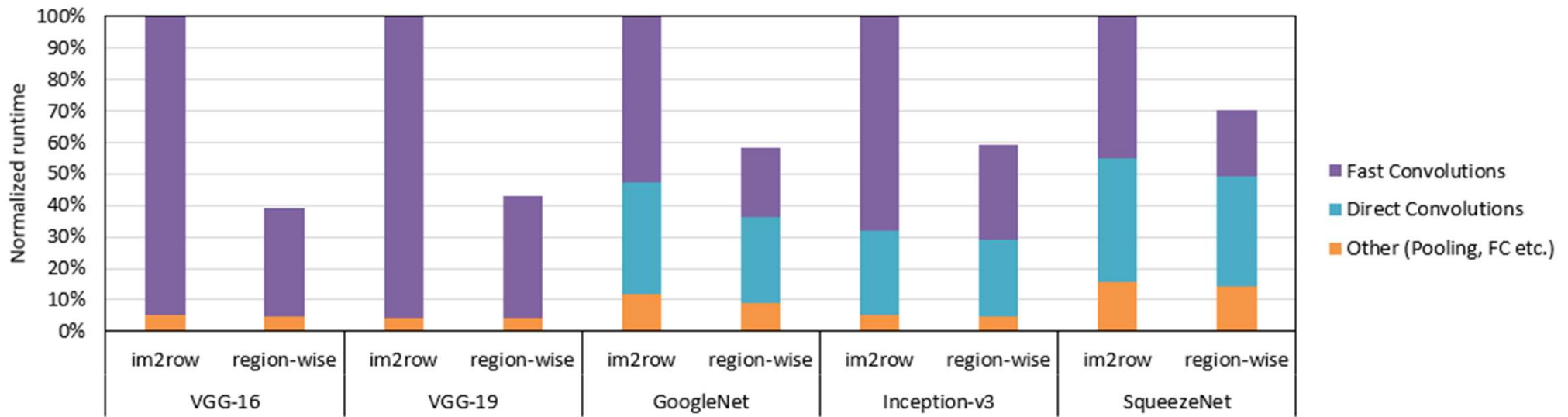
Platform: Huawei HiKey960 Development Platform – 4xA73 cluster

Networks: VGG19, VGG16, GoogleNet, Inception-v3, SqueezeNet

Other: FP32, batchsize 1, 4x multi-threaded through Arm Compute Library (ArmCL)

Measured individual per-layer performance as well as end-to-end run-time, compared with highly optimized conventional ‘im2row GEMM’ convolution strategy

Benchmark Results



Conclusion

- ML is coming to the edge, hard and fast
- ARM CPUs are already widely deployed at the edge, so optimizing for performance here has immediate impact
- Winograd domain is an alternative to conventional im2row/GEMM convolution that reduces math, but requires care to fully realize benefit
- When done properly, can provide as much as a 2.5x speedup on real hardware for end-to-end model inference

Benefits now available in ArmCL!

arm

Thank You

Danke

Merci

谢谢

ありがとう

Gracias

Kiitos

감사합니다

धन्यवाद

شكرًا

תודה



†The Arm trademarks featured in this presentation are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights reserved. All other marks featured may be trademarks of their respective owners.

www.arm.com/company/policies/trademarks