# OpenVX Graph Optimization for Visual Processor Units

**Madushan Abeysinghe, Jesse Villarreal, Lucas Weaver, Jason D. Bakos**

**ASAP 2019**
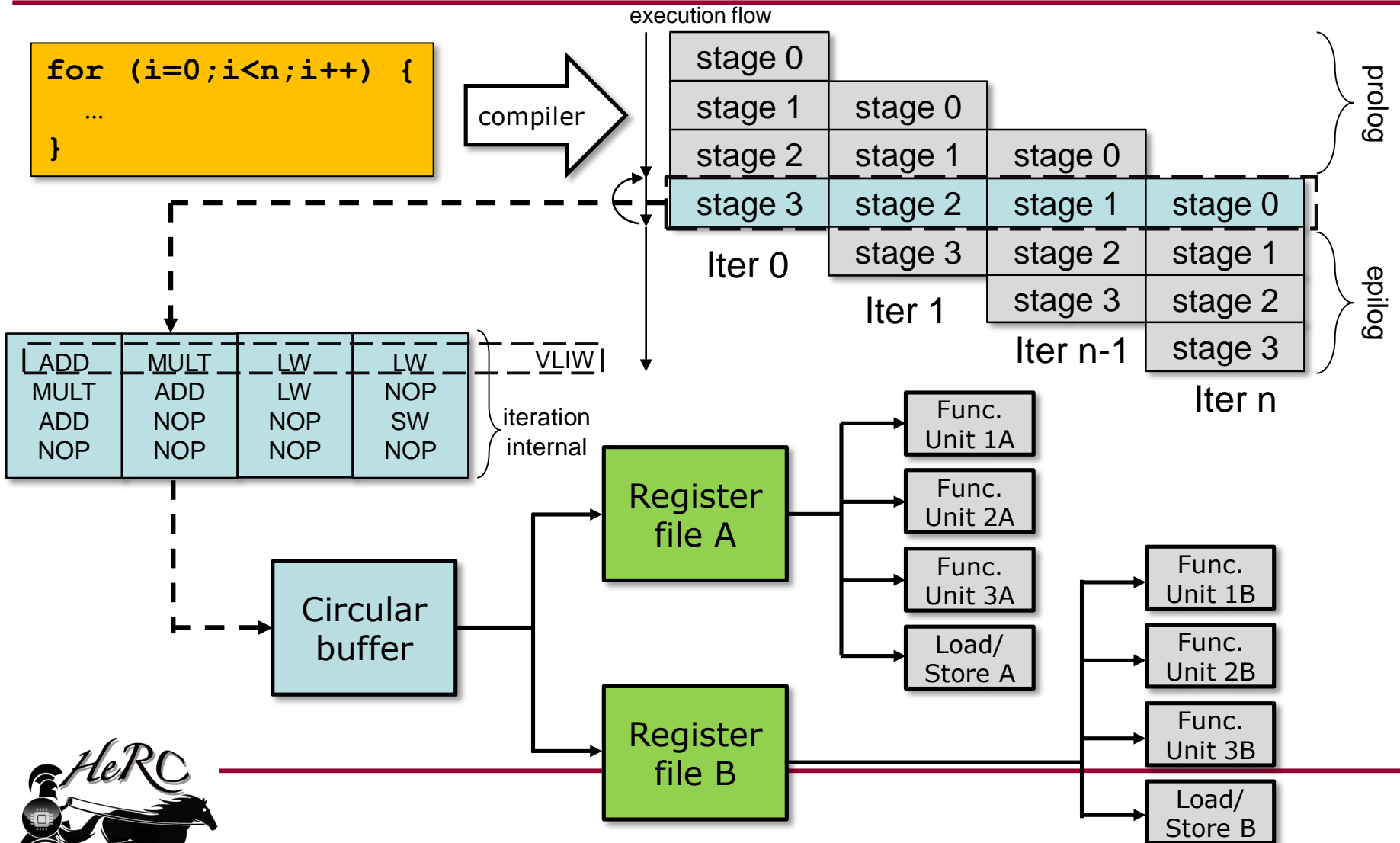**July 16th Cornell Tech, NY**

UNIVERSITY OF
**SOUTH CAROLINA**
1801

*HeRC*

**TEXAS INSTRUMENTS**

Heterogeneous and Reconfigurable
Computing Group

**Presented by Madushan Abeysinghe**

# TI C66 Digital Signal Processor

# OpenVX Framework

- OpenVX is a standardized, cross platform framework to aid in development of accelerated computer vision, machine learning, and other signal processing applications.

- OpenVX allows the programmer to define an application using a graph-based programming model
  - Nodes: highly optimized kernels for the specific architecture
  - Edges: represents data flow

- Code-portable and Performance-portable !

# OpenVX on ADAS TDA2x

# Tiled Processing of OpenVX Graphs

# Tiled Processing of OpenVX Graphs

(out)   (process)   (in)

# DMA and Scratchpad

time →

| input ping buffer | | tile n+1 | tile n+1 |
| input pong buffer | tile n | | tile n+2 | tile n+2 |
| output ping buffer | tile n-1 | | tile n+1 | tile n+1 |
| output pong buffer | | tile n | tile n | tile n+2 |

**processor**

| trigger | process tile n | idle | trigger | process tile n+1 | idle | trigger | process tile n+2 |

**DMA controller**

| idle | write tile n-1 to DRAM | read tile n+1 from DRAM | idle | write tile n to DRAM | read tile n+2 from DRAM | idle | write tile n+1 to DRAM | |

max(trigger+compute,trigger+DMA)

max(compute,DMA)

# Performance Factor 1:  Tile Size



AccumulateSquared Kernel
1600x1200 Input Image

Total Clock Cycles

2.11X

1.58X

32x16  32x24  64x16  64x24  64x48  64x64  128x48  256x24  320x24  384x16

■ DMA cycles  ■ DMA Compute cycles

best compute

best DMA

best
max(DMA,compute)

# Performance Factor 2: Scratchpad Access Pattern



DRAM -> scratchpad (A)

scratchpad -> DRAM (B)

DRAM -> scatchpad (E)

scatchpad -> DRAM (F)

vxColorConvertNode
ex. time = X

vxGaussian3x3Node
ex. time = Z

DRAM -> scratchpad (C)

vxChannelExtractNode
ex. time = Y

scratchpad -> DRAM (D)

time = max(C+D,Y)

scratchpad buffers:  2, size = scratchpad/4

DRAM ->
scratchpad (A)

scratchpad

DRAM ->
scatchpad (E)

scatchpad ->
DRAM (F)

vxColorConvertNode
ex. time = X

vxGaussian3x3Node
ex. time = Z

scratchpad

vxChannelExtractNode
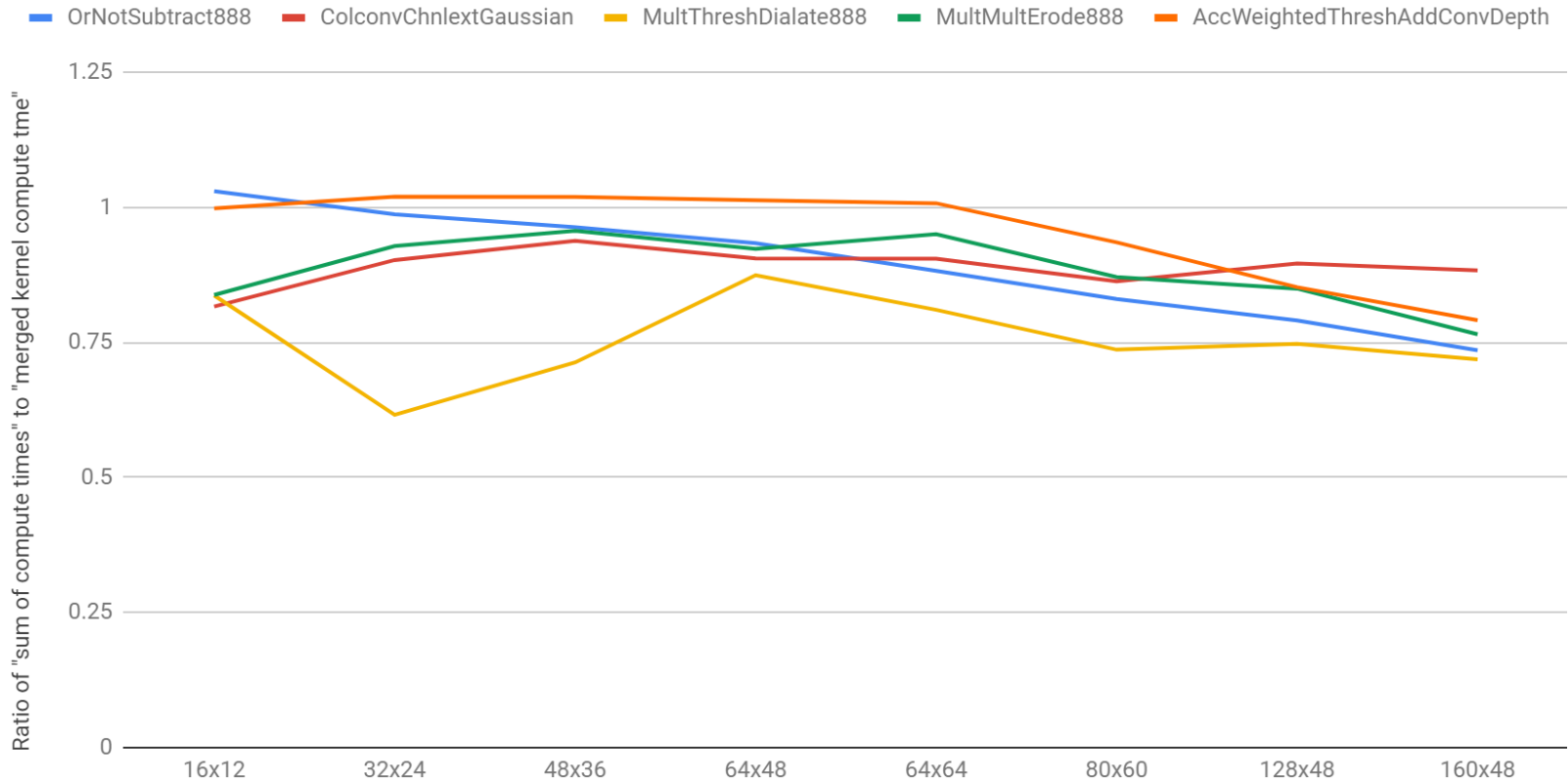ex. time = Y

scratchpad ->
DRAM (D)

time = max(A+D,(X+Y)/λ)
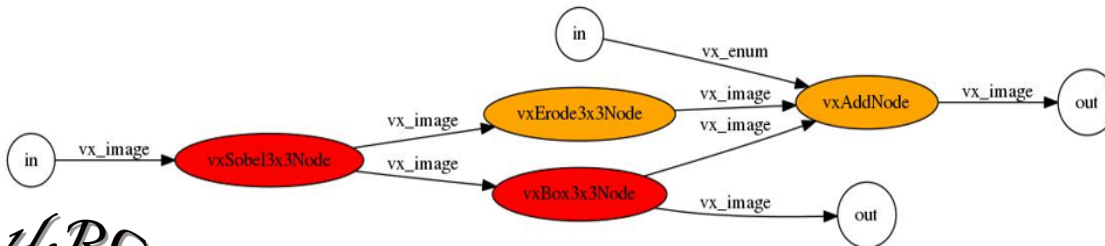
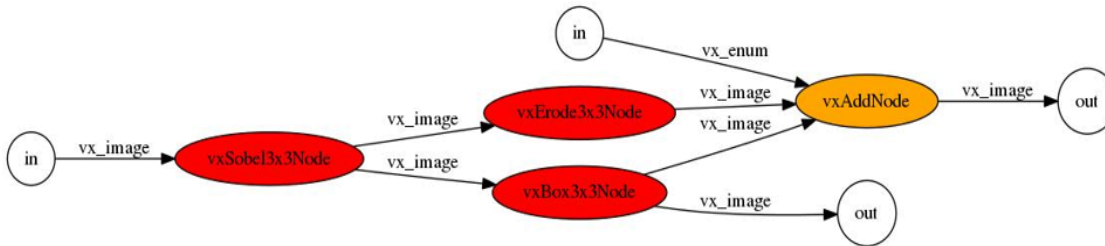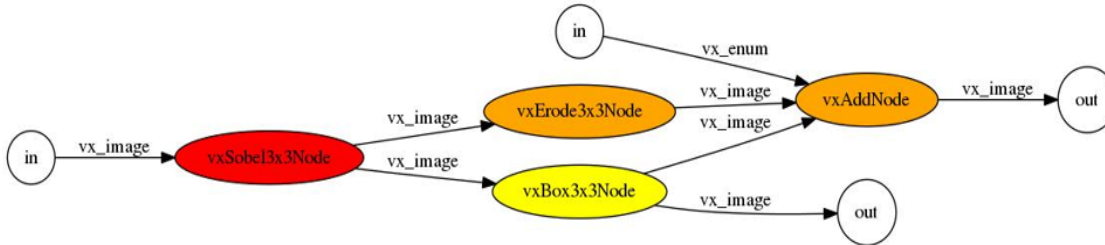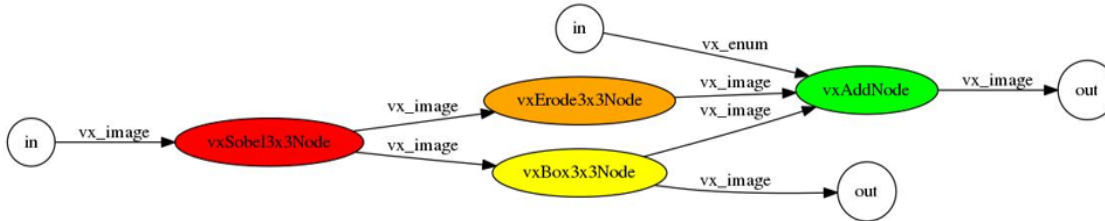scratchpad buffers: 3, size = scratchpad/6

# Lambda

λ Values for Merged Kernels (Observed)

# Graph Partitioning



$$B_n = \sum_{k=0}^{n-1} \binom{n-1}{k} B_k$$

$$B_0 = 1$$

$B_6 = 203$
$B_7 = 877$
$B_8 = 4140$
$B_9 = 21147$
$B_{10} = 115975$
$B_{11} = 678570$
$B_{12} = 4213597$

# Performance Considerations

- Problem:
  - Group size determines:
    - Number of DMA transactions
    - Maximum tile size
    - Lambda
  - Tile size determines:
    - DMA bandwidth
    - Compute throughput

- Approach:
  Train memory performance model and compute model to optimize grouping and tile size together

# DMA/Memory System Performance Model

- Characterize:
  - 19 kernels (originally 35), 25 tile sizes, all possible pixel depths (per kernel)
  - **1460** total test cases (culled from 3240 original tests)
- Features:
  1. input tile width, height in pixels (**TIW, TIH**)
  2. output tile width, height in pixels (**TOW**, **TOH)**
  3. total input, output width (tile width x pixel size x number of inputs) (**TTIW, TTOW**)
  4. total input, output size, TTIW*TIH (**TIS, TOS**)
  5. number of inputs, outputs (**NI, NO**)
  6. pixel depth (**PD**)
  7. stencil neighborhood width, height (**SNW, SNH**)

- Best accuracy: **TTIW, TIH, TTOW, TOH**

# Performance Models

- **DMA: linear interpolant model of DMA b/w**
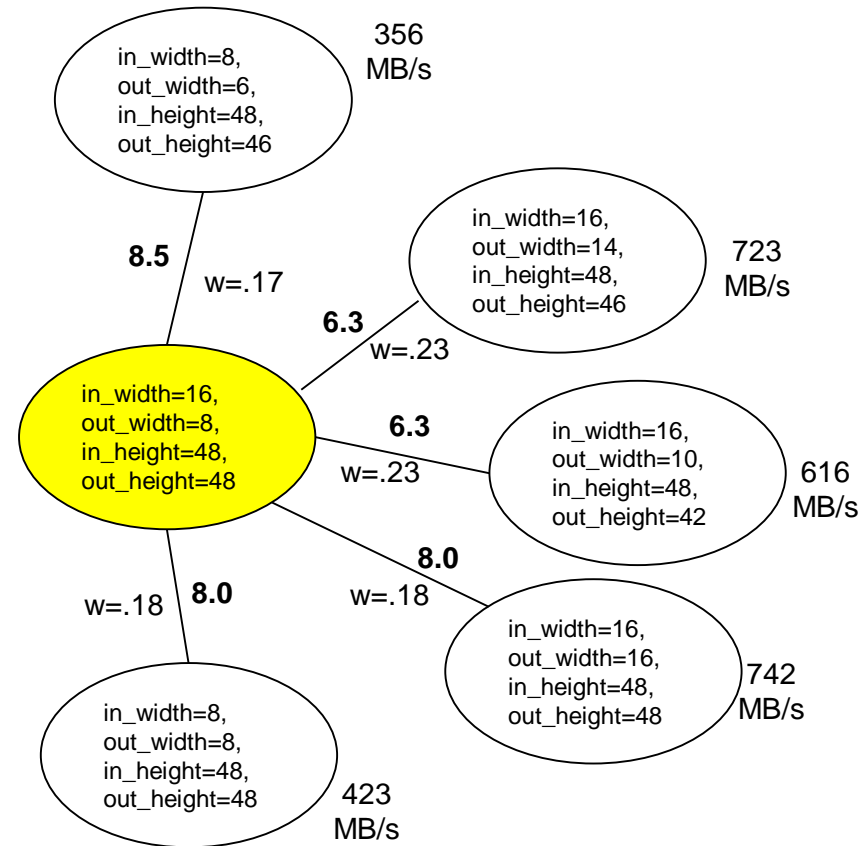  - Find nearest 5 observations
    - Calculate weights
      - $w_i = \dfrac{\frac{1}{distance_i}}{\sum_j \frac{1}{distance_j}}$
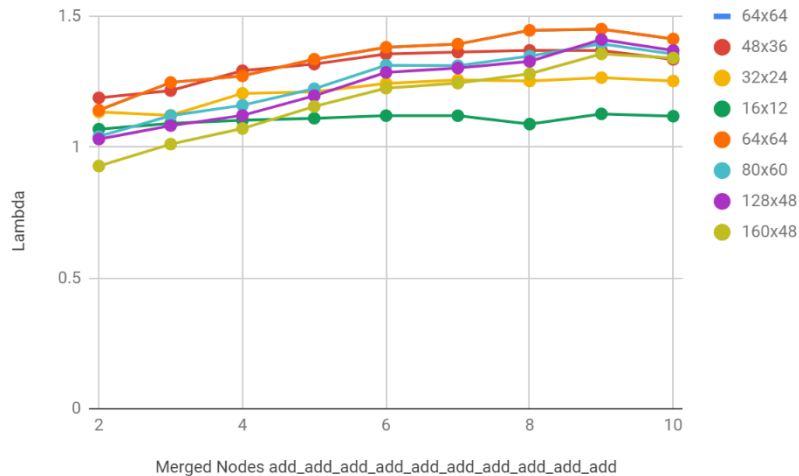    - Calculate weighted average

  - RMS training error = 32 MB/s
  - RMS testing error = 248 MB/s
  - Range = 189 MB/s to 4.33 GB/s

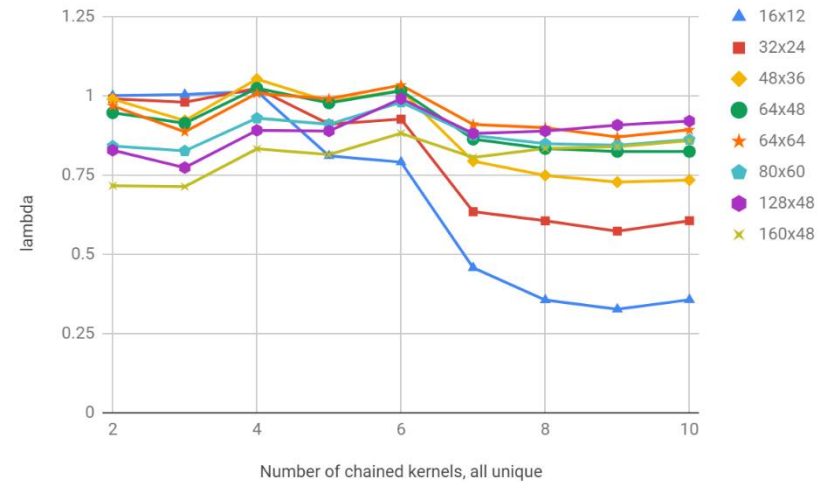- **Compute: lookup kernel names and nearest tile WxH prod / Lambda**

in_width=8,
out_width=6,
in_height=48,
out_height=46
356 MB/s

in_width=16,
out_width=14,
in_height=48,
out_height=46
723 MB/s

in_width=16,
out_width=8,
in_height=48,
out_height=48

in_width=16,
out_width=10,
in_height=48,
out_height=42
616 MB/s

in_width=16,
out_width=16,
in_height=48,
out_height=48
742 MB/s

in_width=8,
out_width=8,
in_height=48,
out_height=48
423 MB/s

**8.5** w=.17

**6.3** w=.23

**6.3** w=.23

**8.0** w=.18

**8.0** w=.18

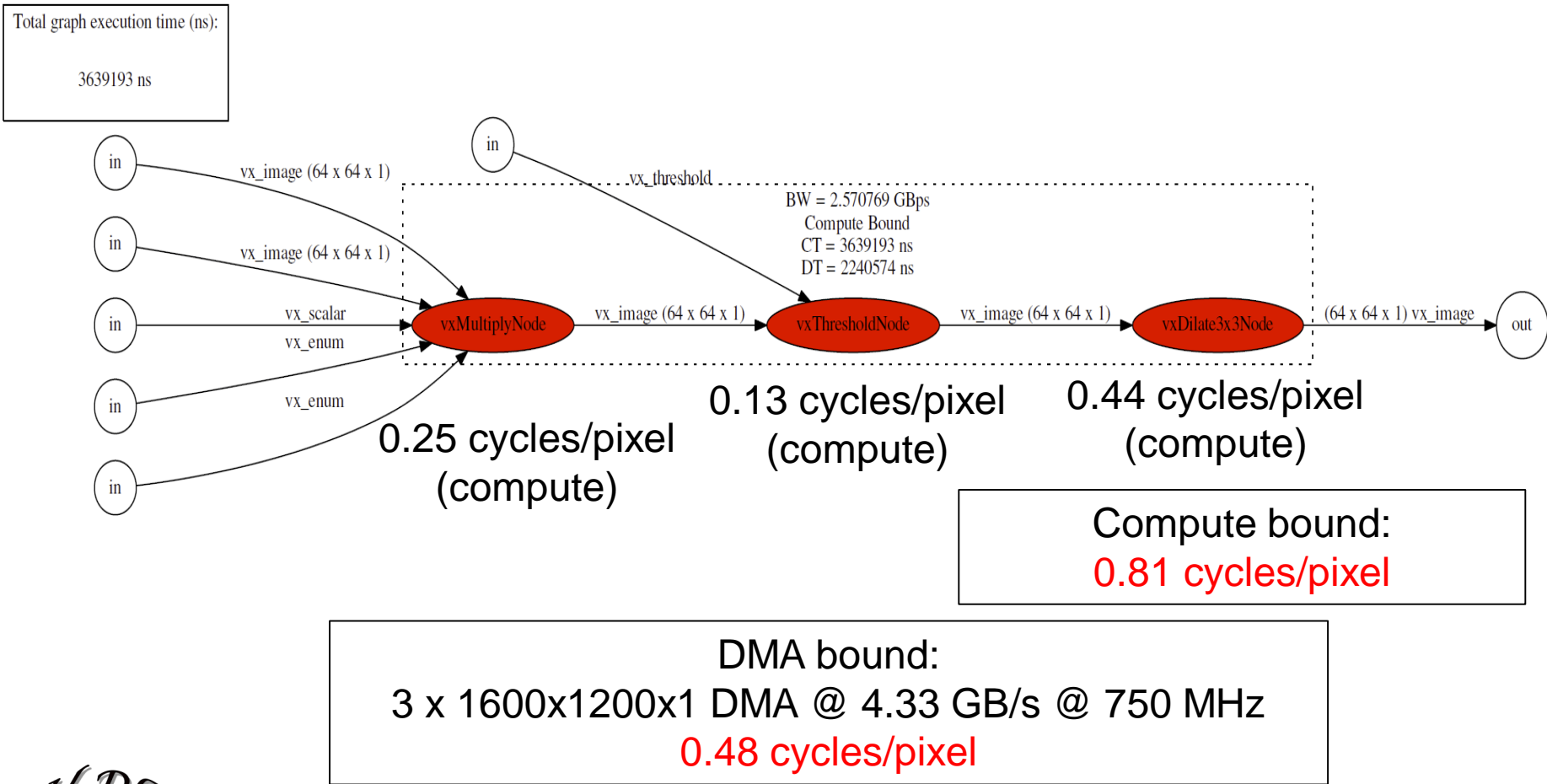# Lamdba Model

Lamdba vs graph size: same kernel
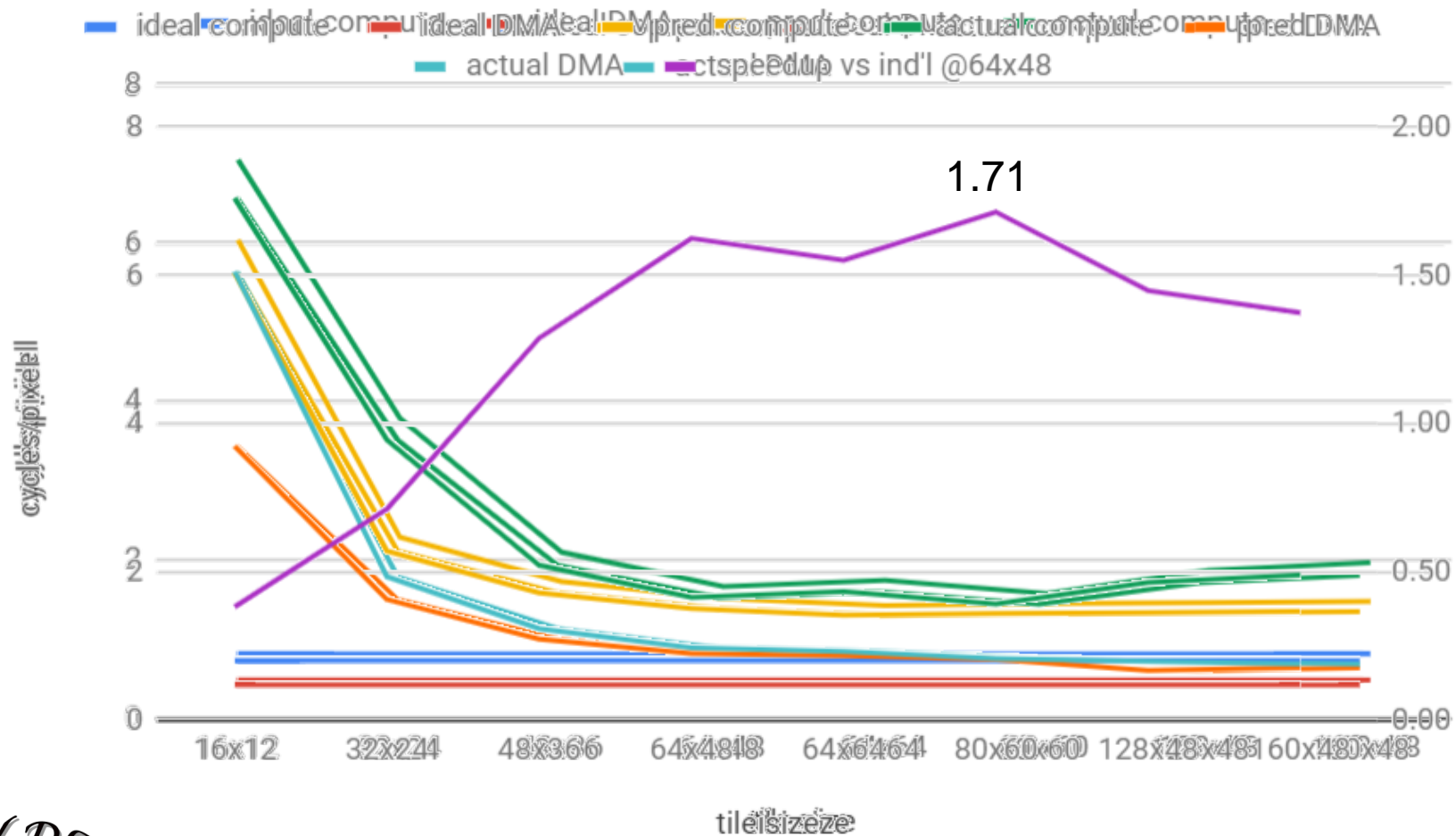


Lamdba vs graph size: unique kernels



- Lambda value depends on number of kernels and number of unique kernels
- Used a piecewise interpolated model
  - RMS error = 0.03 for training set of 43 test graphs

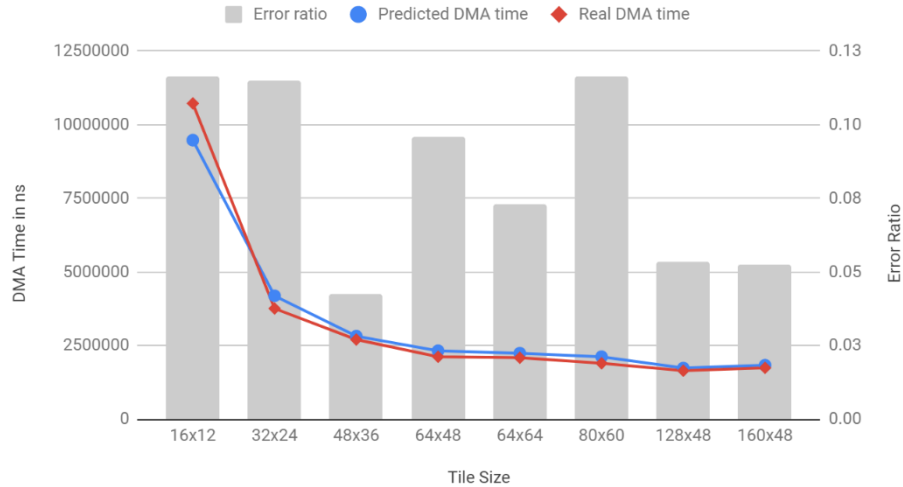# Model Accuracy and Performance



Total graph execution time (ns):

3639193 ns

in — vx_image (64 x 64 x 1)

in — vx_image (64 x 64 x 1)

in — vx_scalar / vx_enum

in — vx_enum

in

vx_threshold

BW = 2.570769 GBps
Compute Bound
CT = 3639193 ns
DT = 2240574 ns

vxMultiplyNode — vx_image (64 x 64 x 1) — vxThresholdNode — vx_image (64 x 64 x 1) — vxDilate3x3Node — (64 x 64 x 1) vx_image — out

0.25 cycles/pixel
(compute)

0.13 cycles/pixel
(compute)

0.44 cycles/pixel
(compute)

Compute bound:
0.81 cycles/pixel

DMA bound:
3 x 1600x1200x1 DMA @ 4.33 GB/s @ 750 MHz
0.48 cycles/pixel

# Impact of Tile Size on Merged Kernels
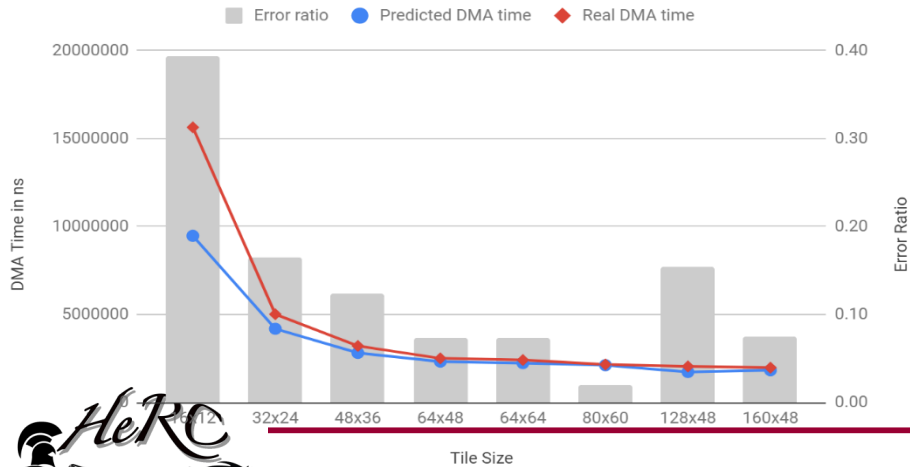
# Model Accuracy



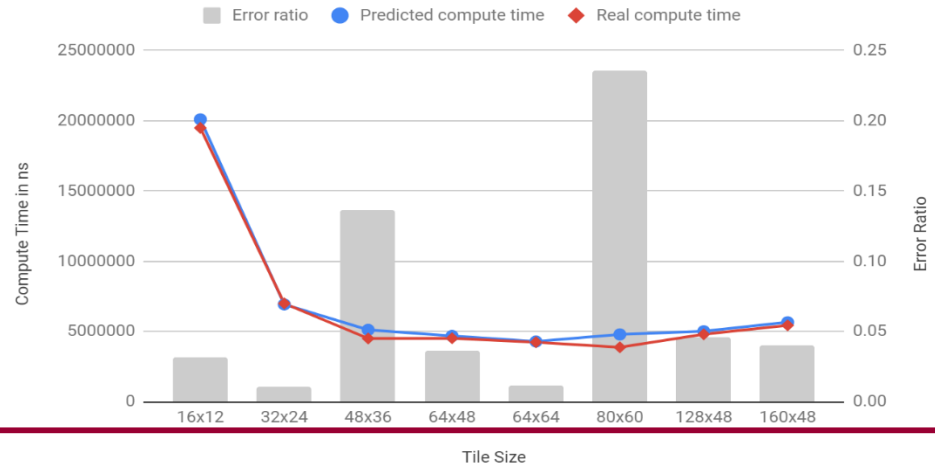DMA time comparison for Or_Not_Subtract Merged Kernel

Compute time comparison for Or_Not_Subtract Merged Kernel

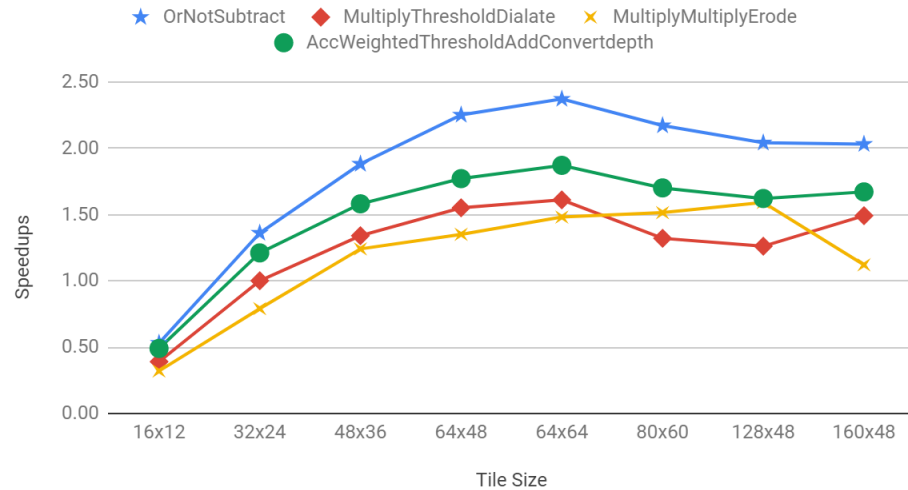DMA time comparison for Multiply_Multiply_Erode Merged Kernel

Compute time comparison for Multiply_Multiply_Erode Merged Kernel
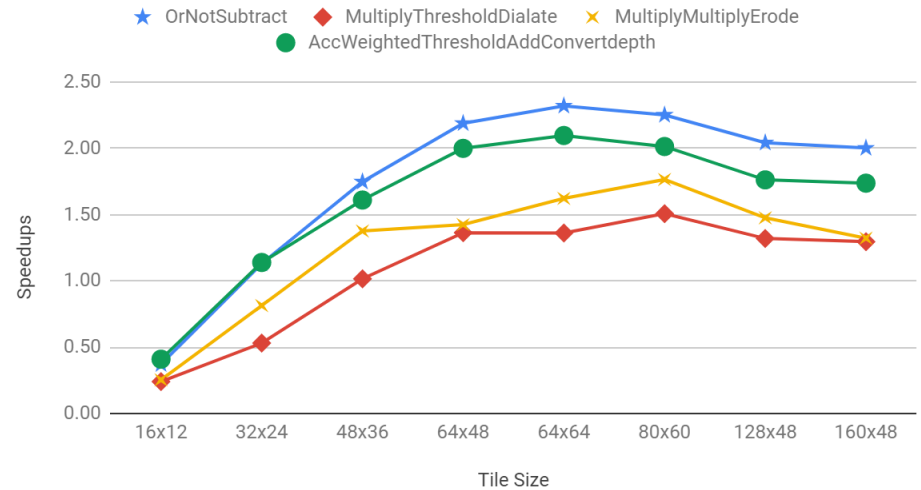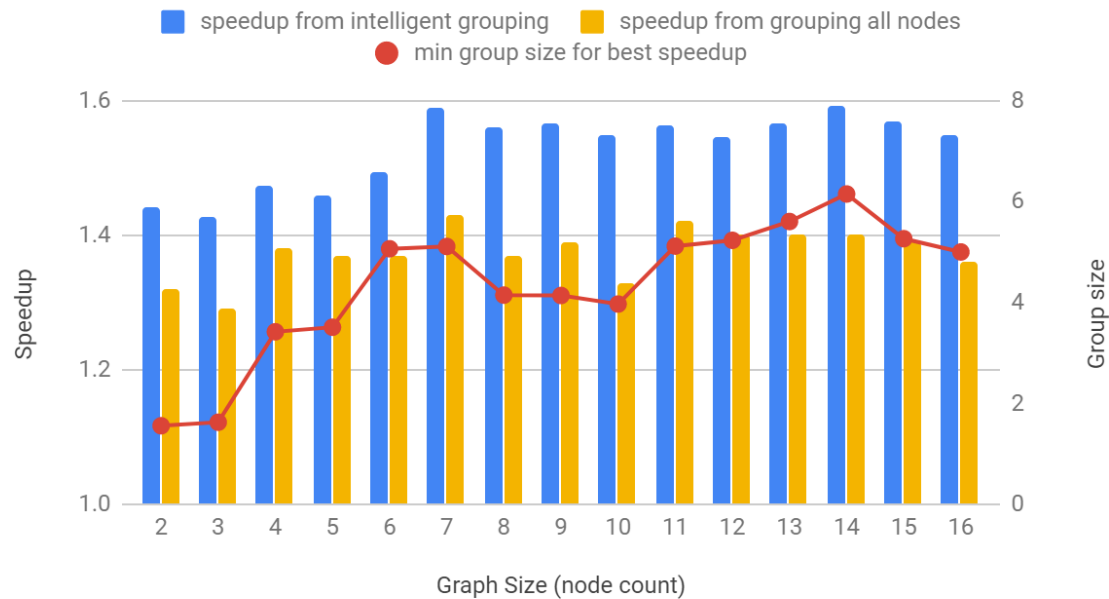
# Speedups



Comparison of Predicted and Real Speedups Speedups
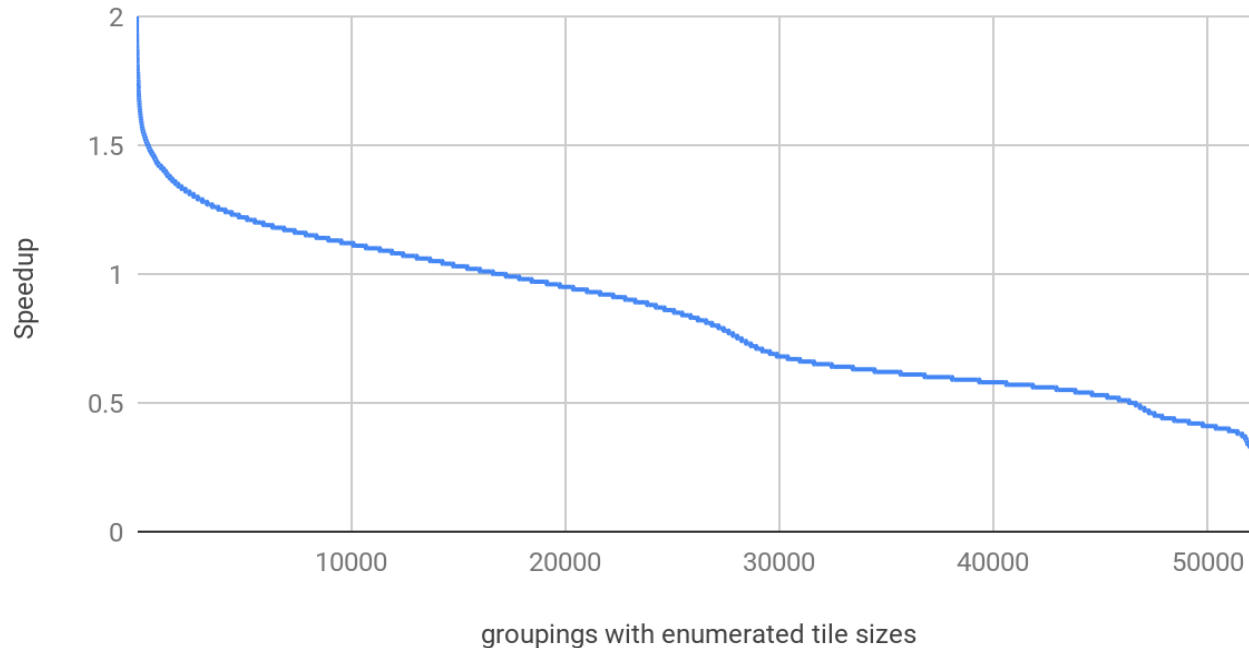
# Randomly Generated Graphs

- For each graph
  - For each grouping
    - For each allocable tile size (constrained by scratchpad capacity)…
      - Find predicted speedup (vs non merged kernels with 64x48 tiles)

Speedup vs. Graph Size

# Speedup Distribution

5 Node Graph (Predicted Speedup = 2.00)



groupings with enumerated tile sizes

- Speedups above 1.5 = 0.95%
- Speedups between 1 and 1.5 = 30.74%
- Speedups below 1 = 68.30%

# Conclusions

- Performance models for DMA time and compute time for OpenVX graphs on C66 DSP
  - Accurate to within 13% for DMA model and 8% for compute model.
- When used to automatically group random graphs, achieves a speedup of ~1.6 on average

- Future Work:
  - Develop models to predict performance from merging the loops for each OpenVX kernel

# Thank you!

Questions?