

Accelerating Pooling through Im2col and Col2im Instructions in the DaVinci Architecture

**Caio S.
Rohwedder**
Unicamp - Brazil

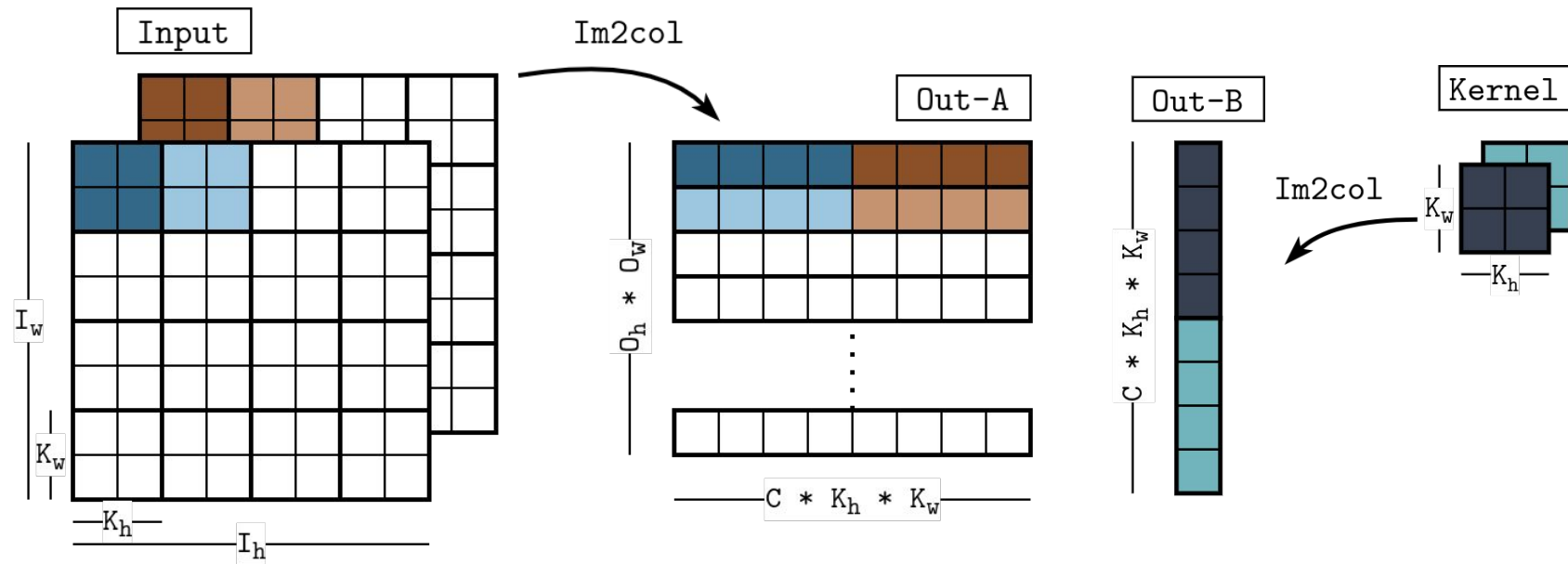
José N. Amaral
University of
Alberta

Guido Araújo
Unicamp - Brazil

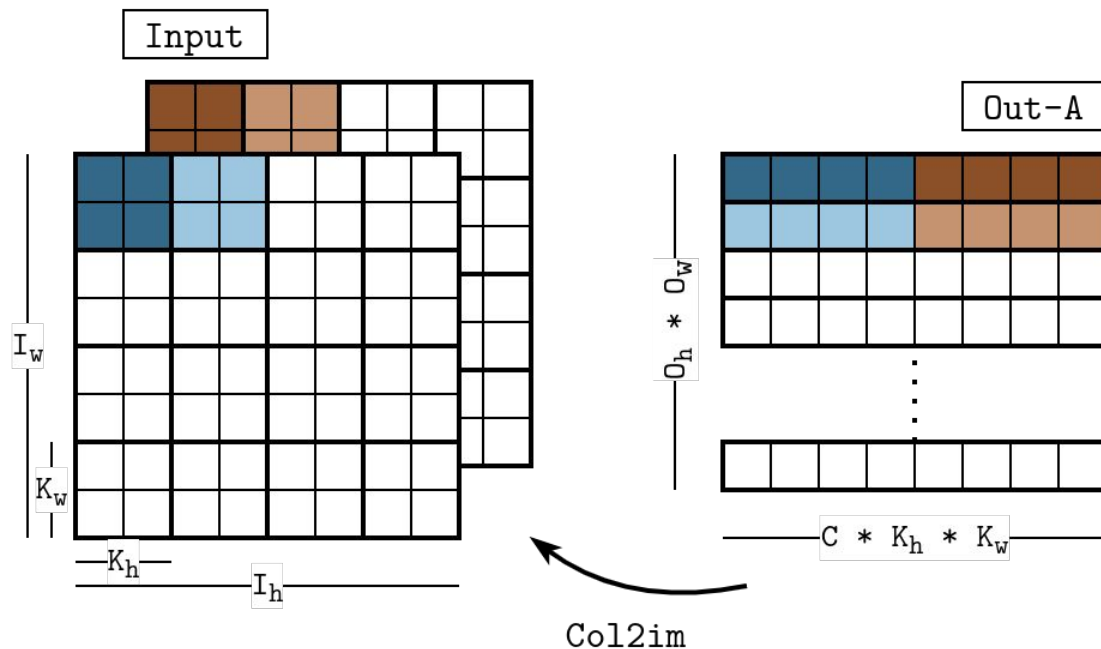
Amy Wang
Toronto
Heterogeneous
Compiler Lab

**Giancarlo
Colmenares**
Toronto
Heterogeneous
Compiler Lab

Im2col + GEMM

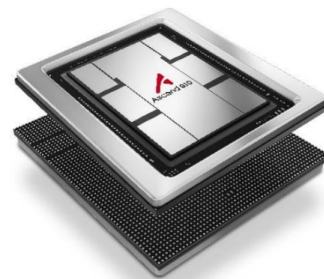


Col2im: Backward Operator



DL Accelerators

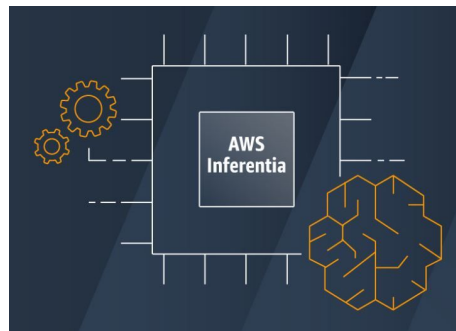
Huawei's Ascend 910²



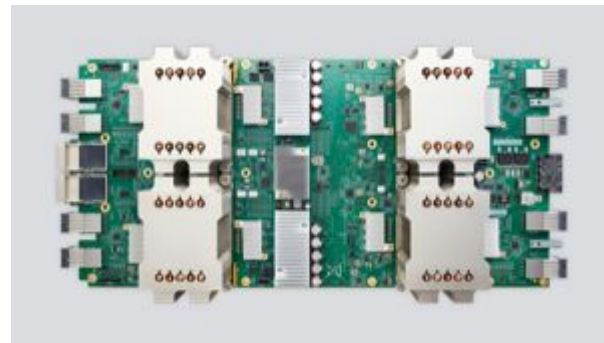
Apple's A12³



Amazon's Inferentia⁴



Google's TPU v2¹

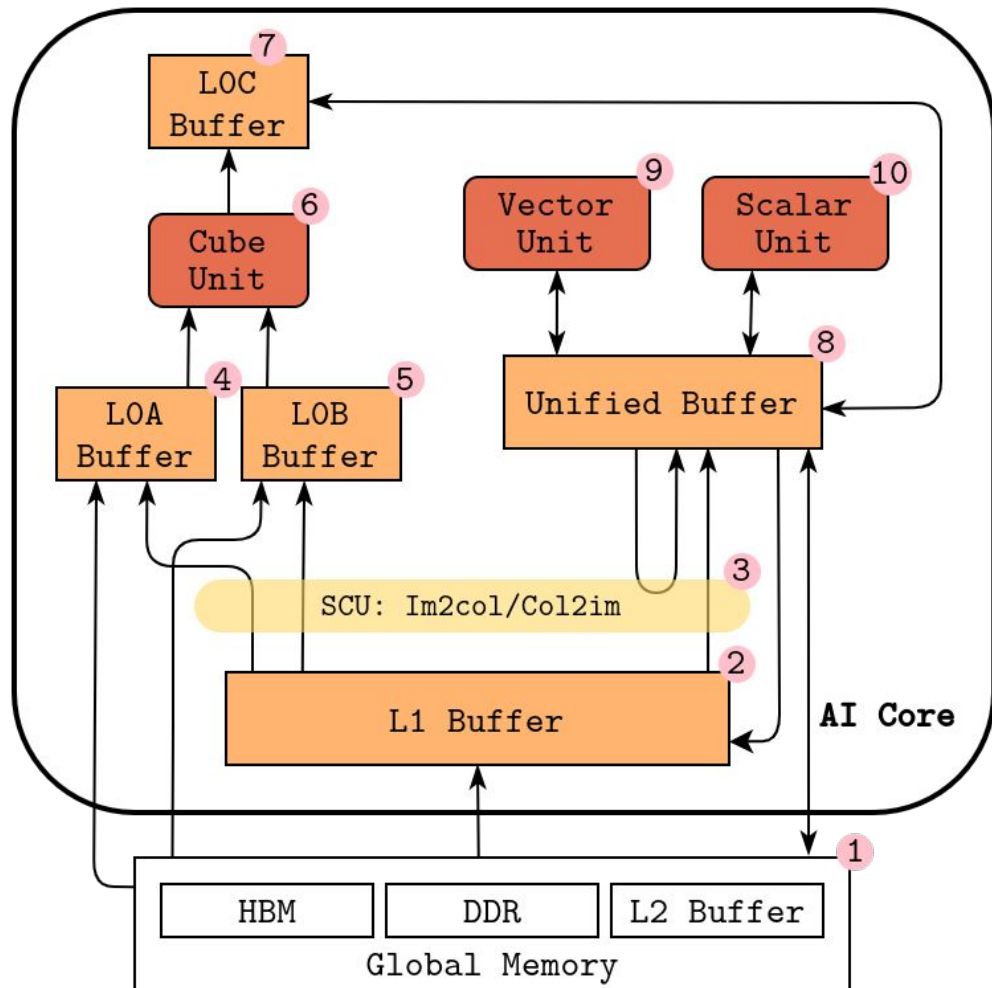


¹<https://cloud.google.com/tpu> ² <https://e.huawei.com/en/products/cloud-computing-dc/atlas/ascend-910>

³<https://www.techinsights.com/blog/apple-iphone-xs-max-teardown> ⁴<https://aws.amazon.com/machine-learning/inferentia/>

DaVinci - Ai Core

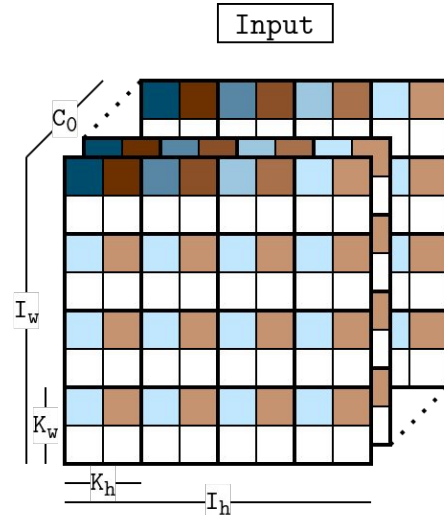
$NCHW \rightarrow NC_1HWC_0$



Im2col Instruction

— — —

$$(x, y) = (0, 0)$$

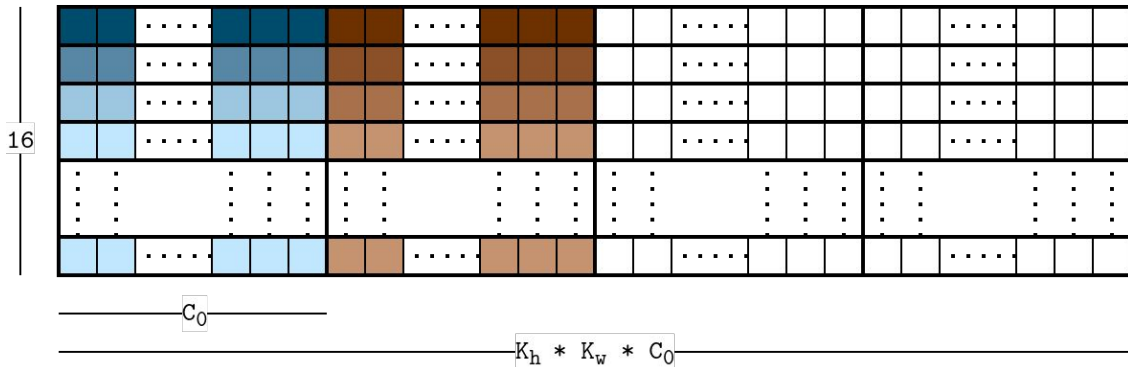


Im2col_Load
 $(x_k, y_k) = (0, 0)$

Im2col_Load
 $(x_k, y_k) = (0, 1)$

Im2col_Load
 $(x_k, y_k) = (1, 0)$

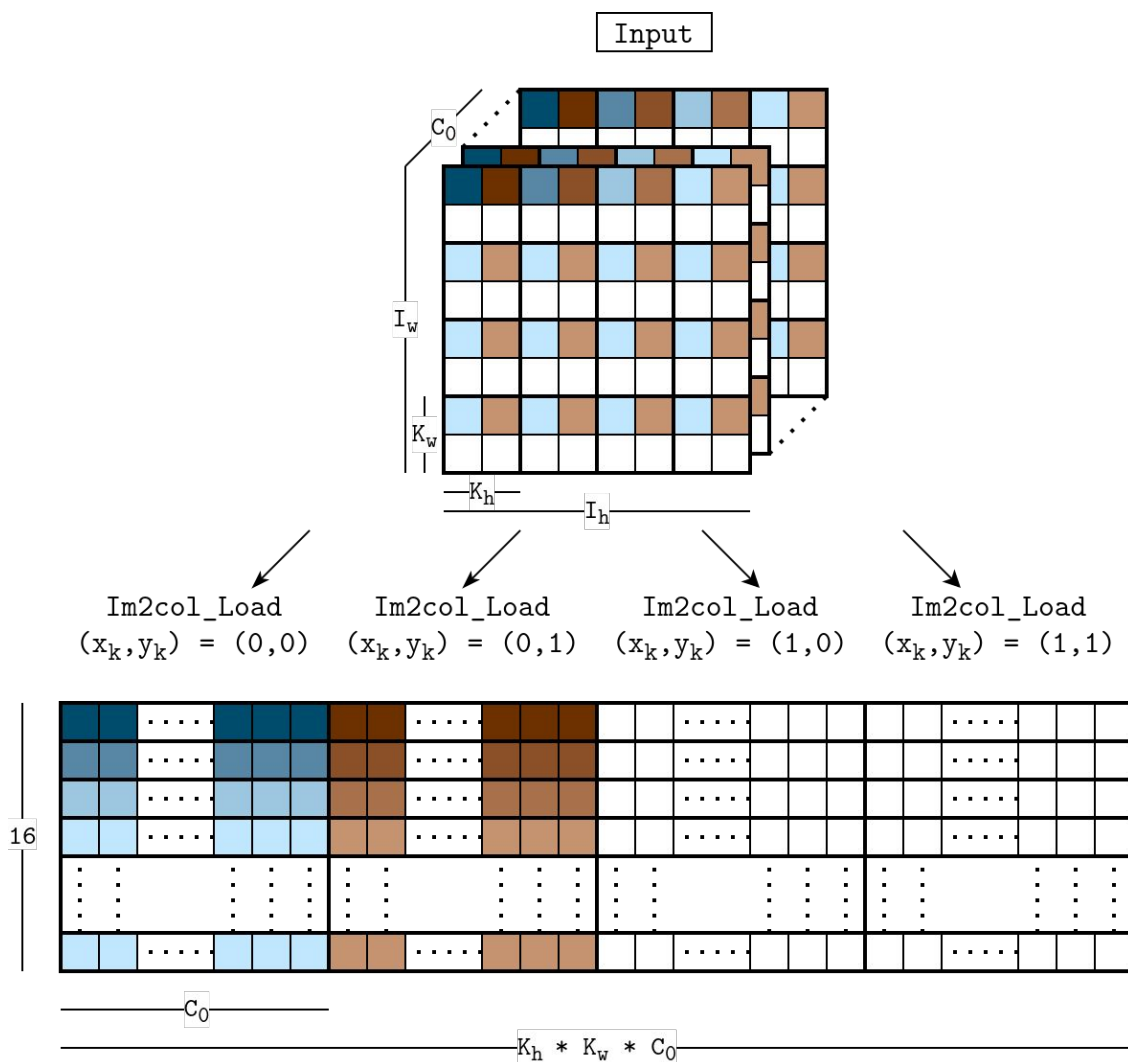
Im2col_Load
 $(x_k, y_k) = (1, 1)$



Im2col Instruction

$(x, y) = (0, 0)$

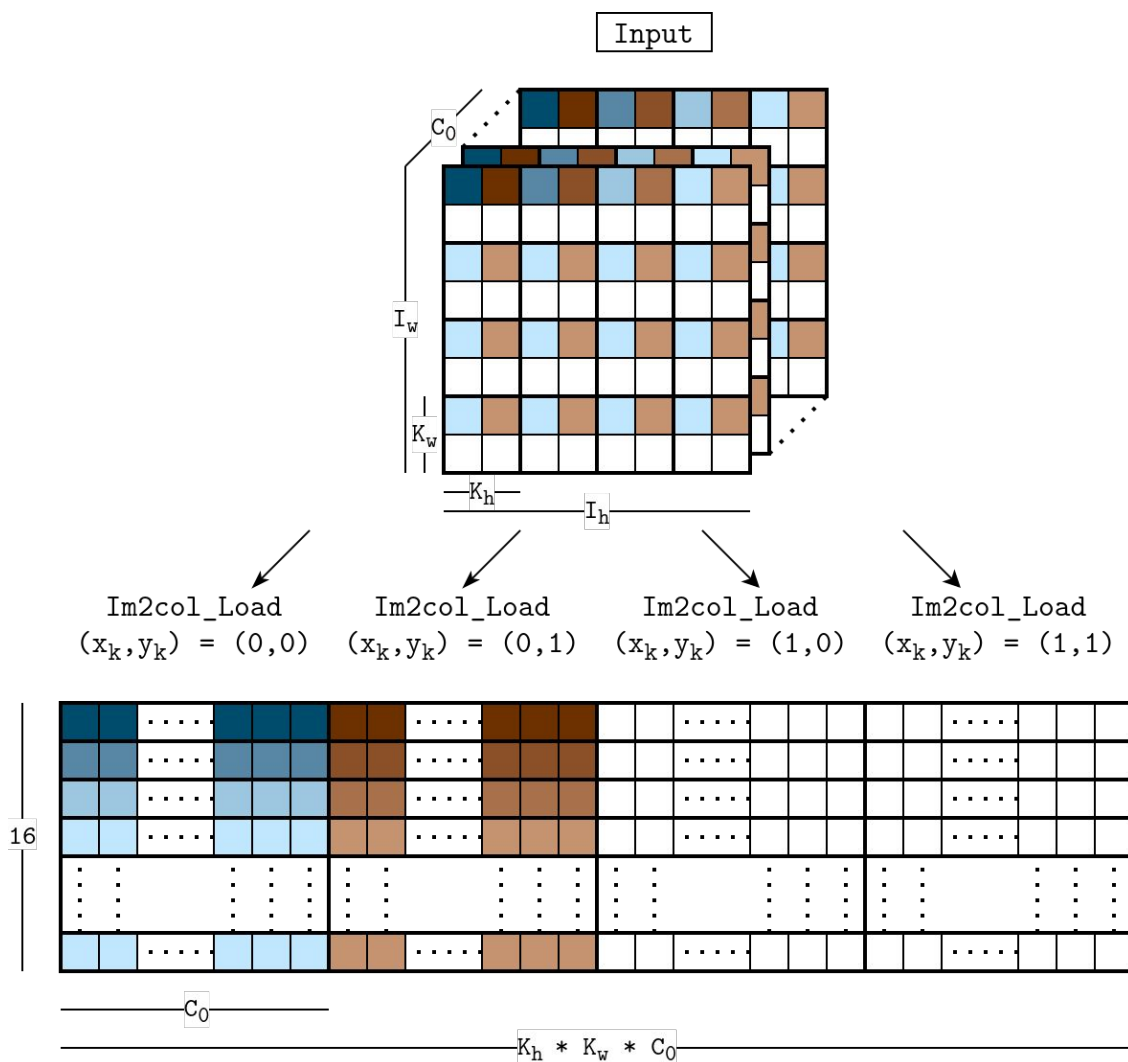
1. Select the next 16 patches from (x, y)



Im2col Instruction

$(x, y) = (0, 0)$

1. Select the next 16 patches from (x, y)
2. Select position (x_k, y_k) relative to the patches

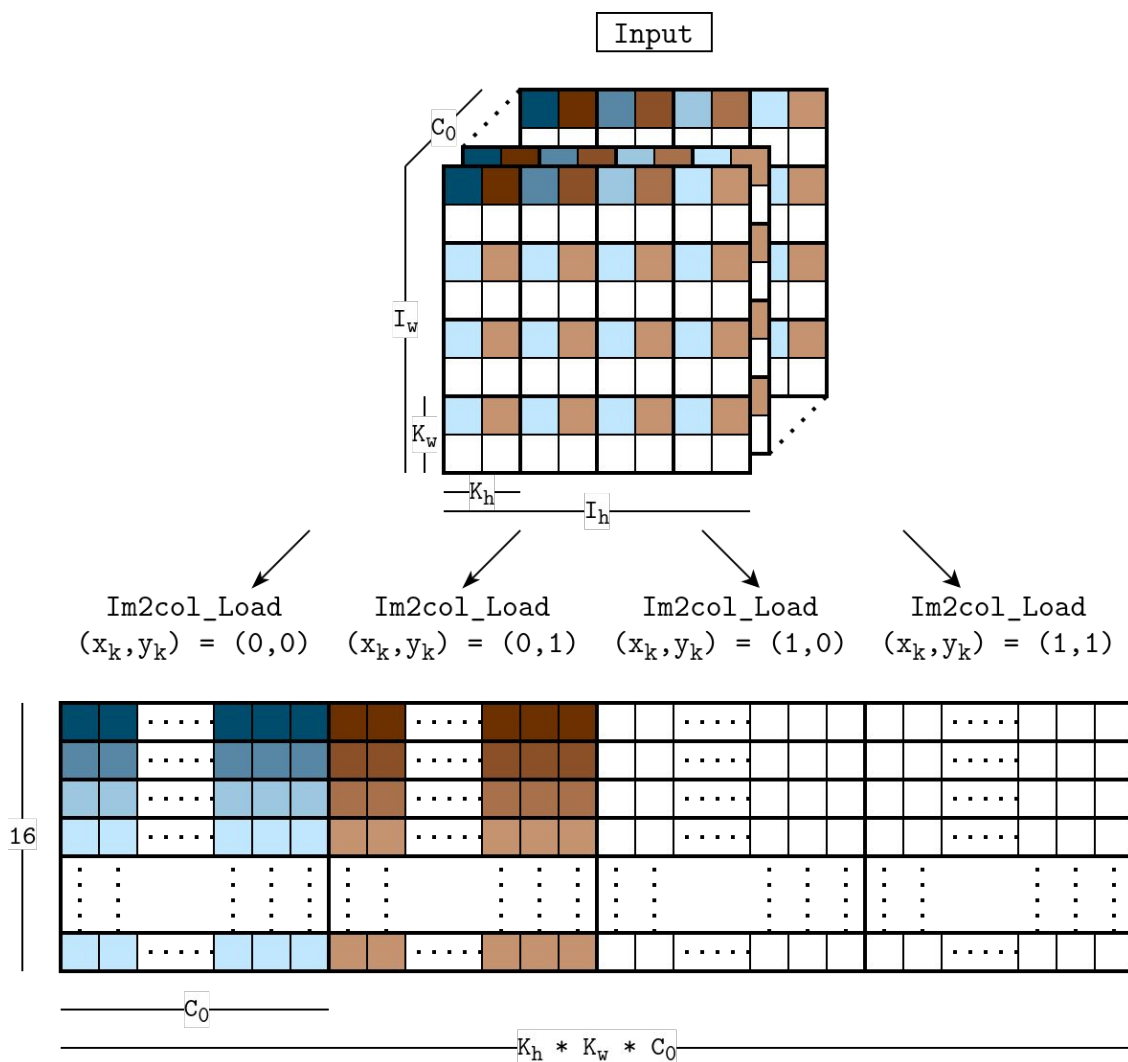


Im2col Instruction

— — —

$(x, y) = (0, 0)$

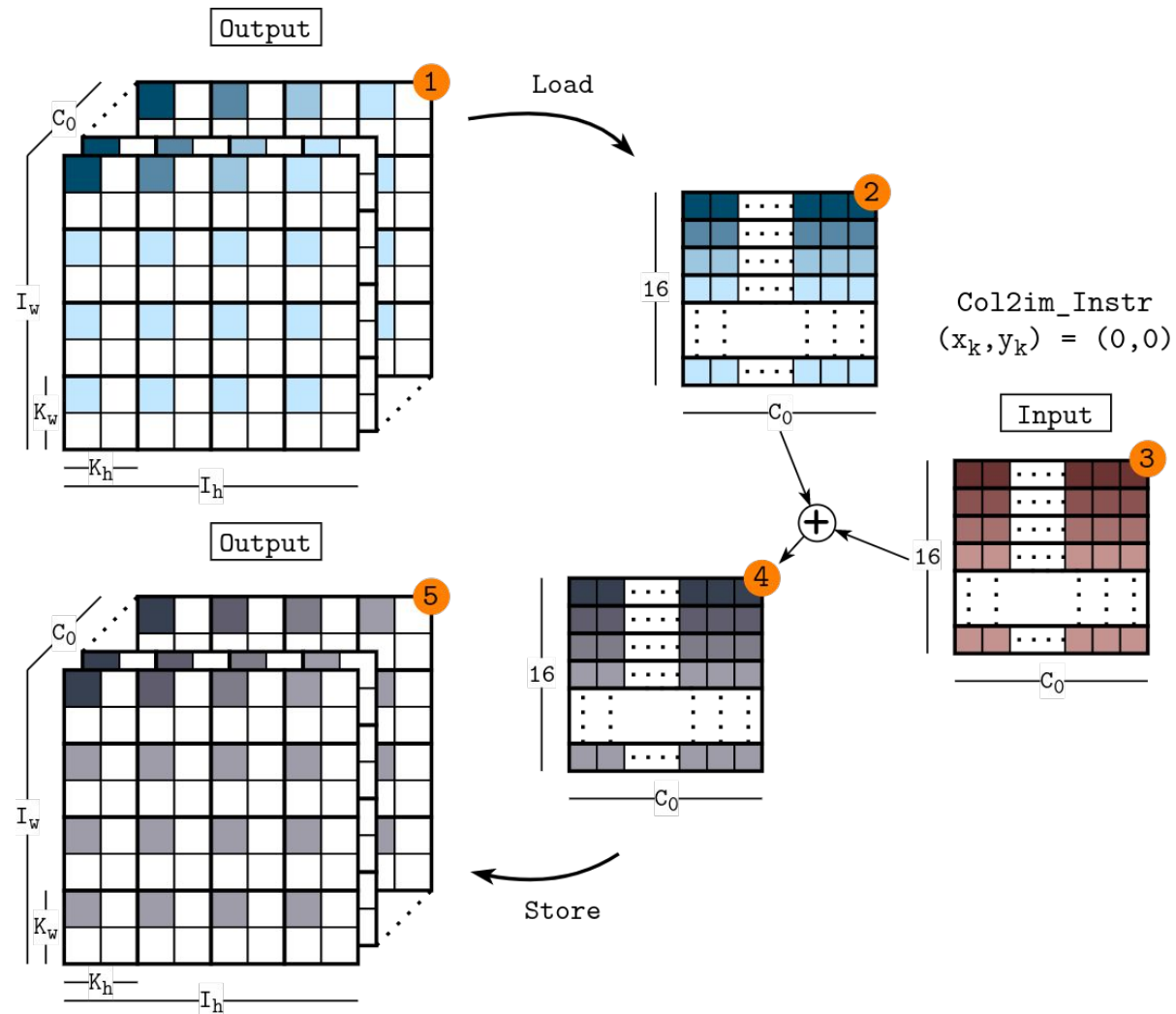
1. Select the next 16 patches from (x, y)
2. Select position (x_k, y_k) relative to the patches
3. Load all C_0 channels from the selected positions



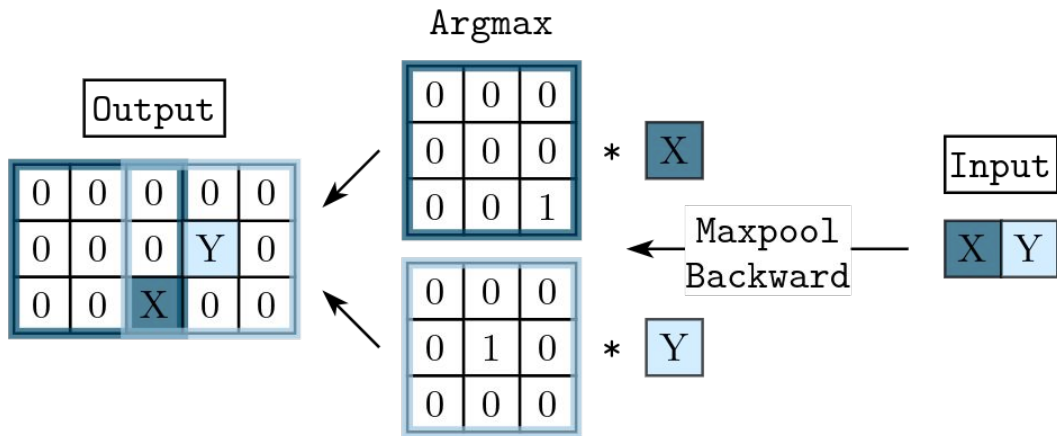
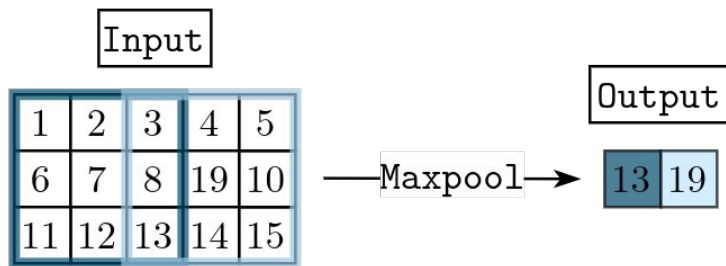
Col2im Instruction

Initialize output
with zeros

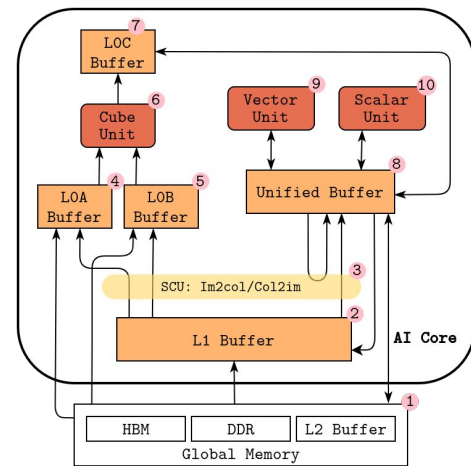
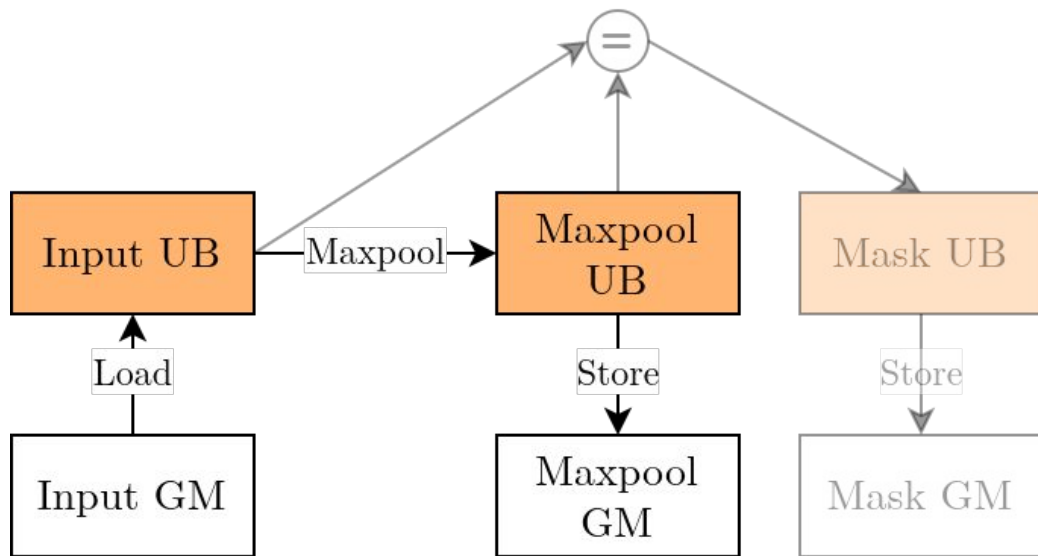
$(x, y) = (0, 0)$



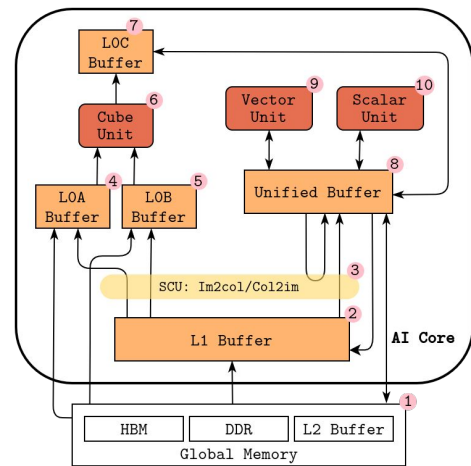
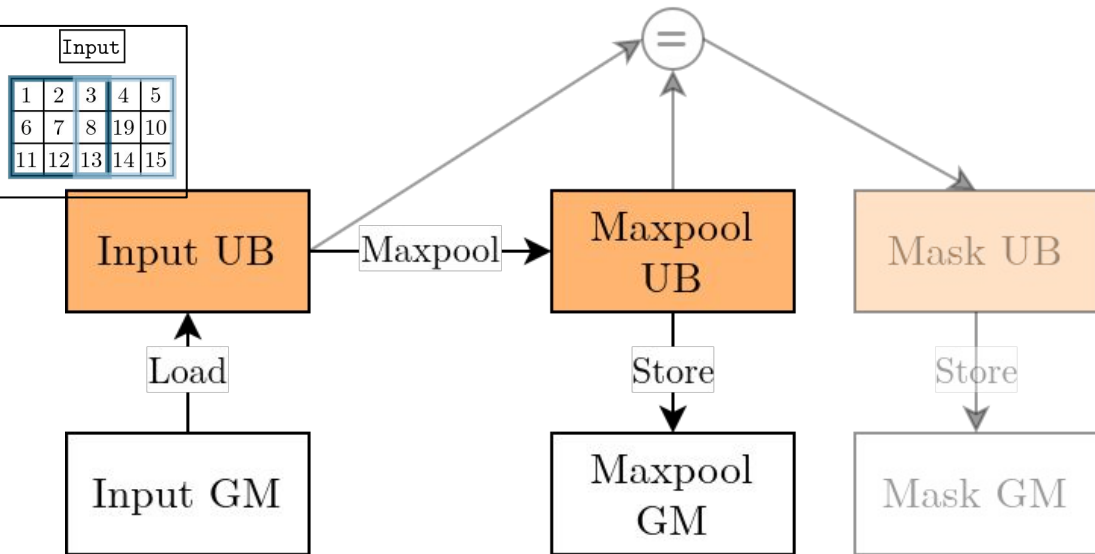
Pooling Operators



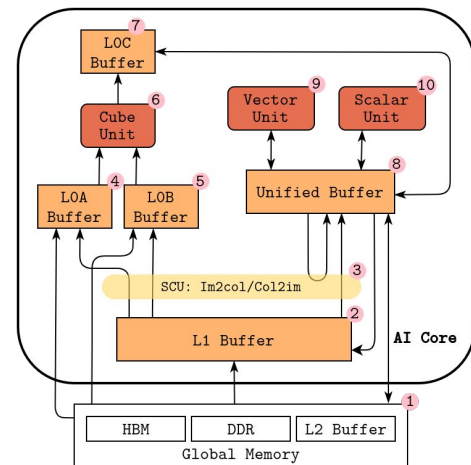
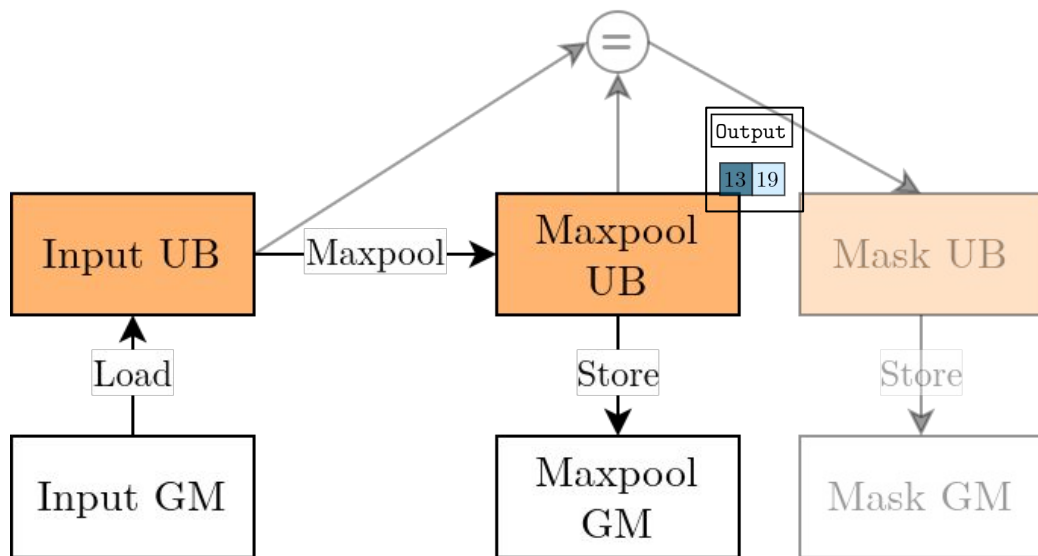
Pooling for DaVinci



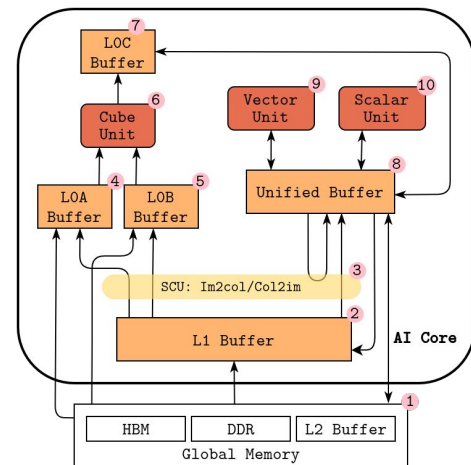
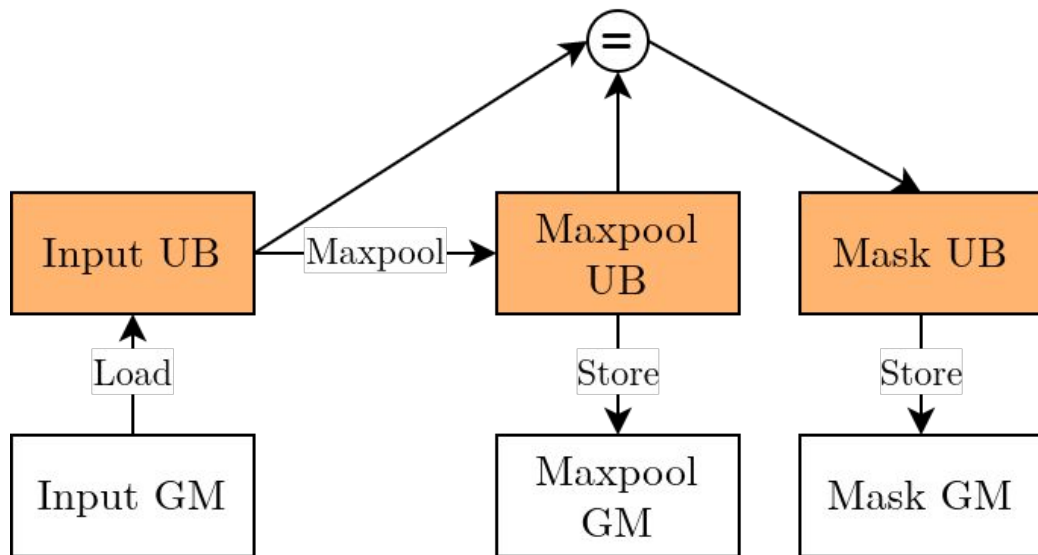
Pooling for DaVinci



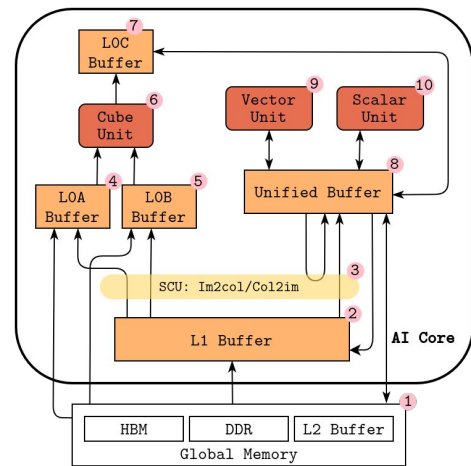
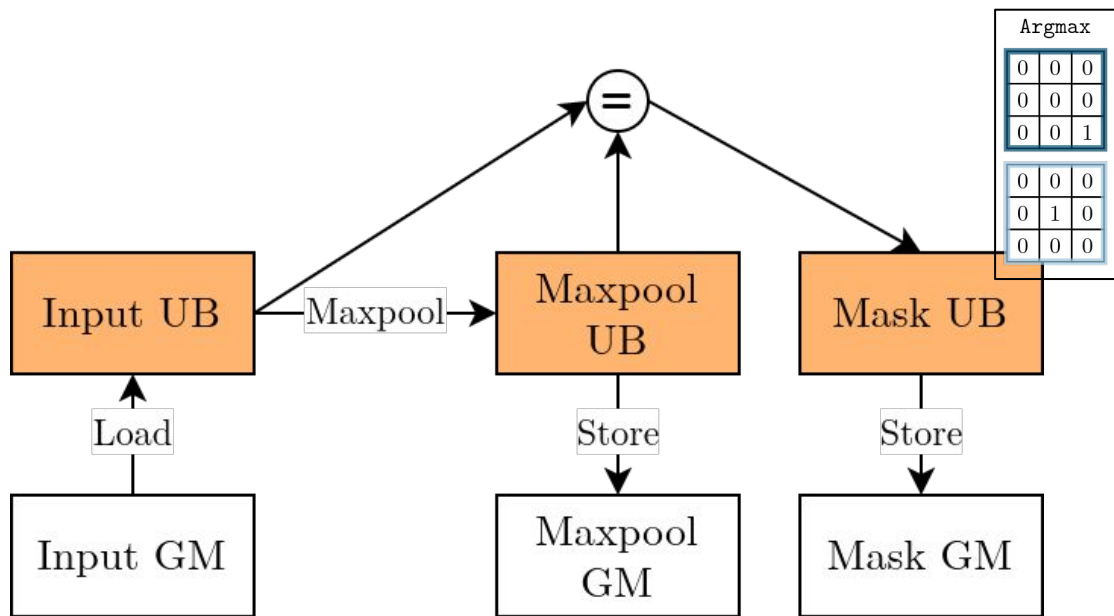
Pooling for DaVinci



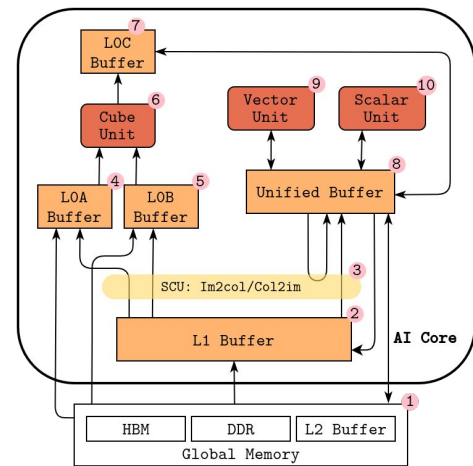
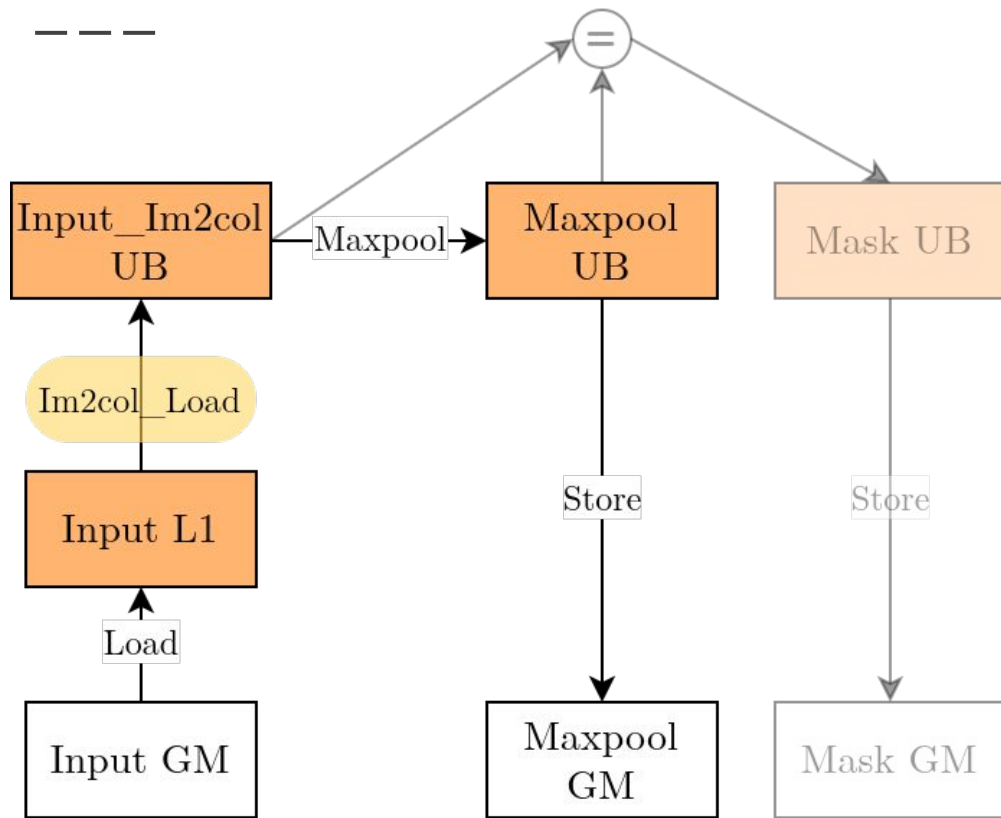
Pooling for DaVinci



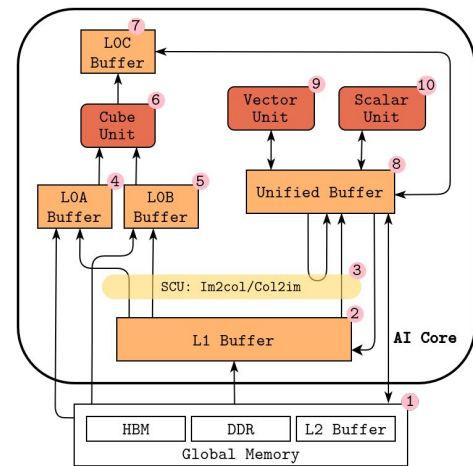
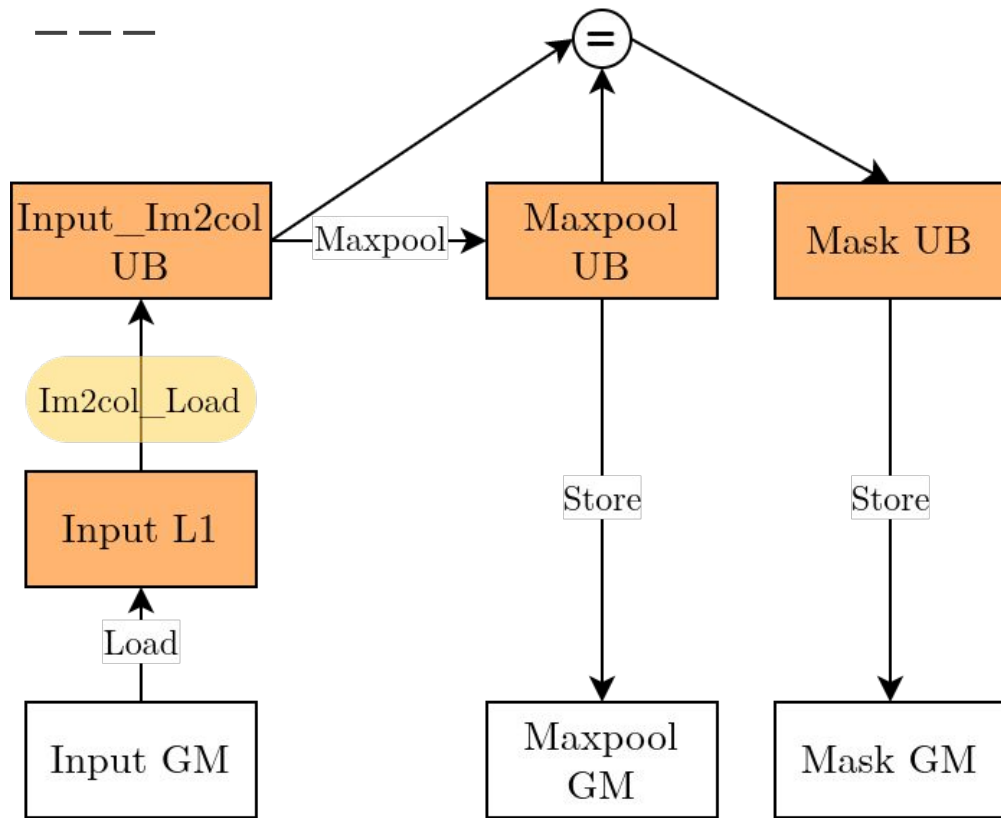
Pooling for DaVinci



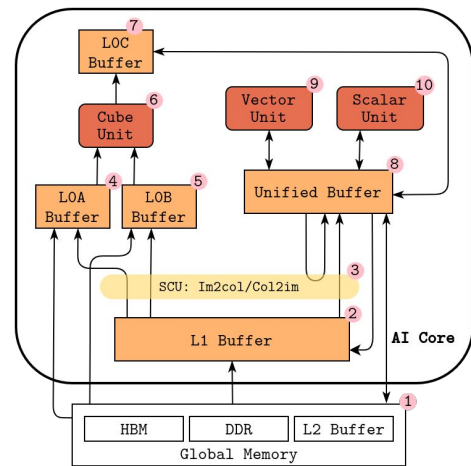
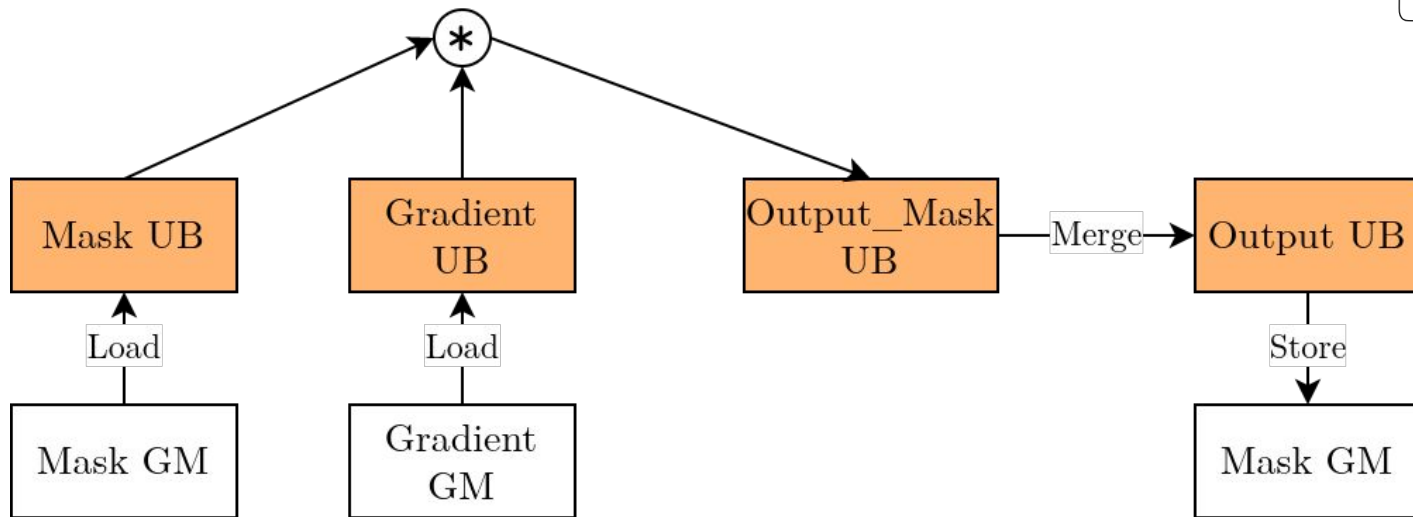
Im2col Based Pooling



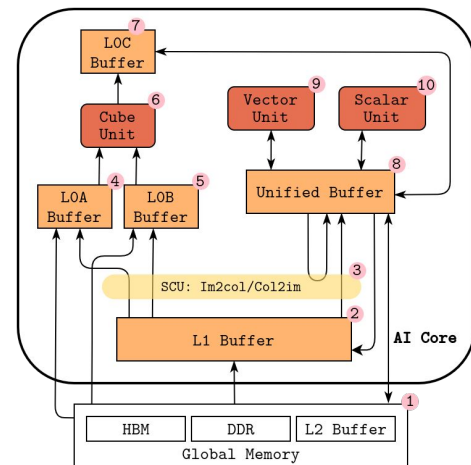
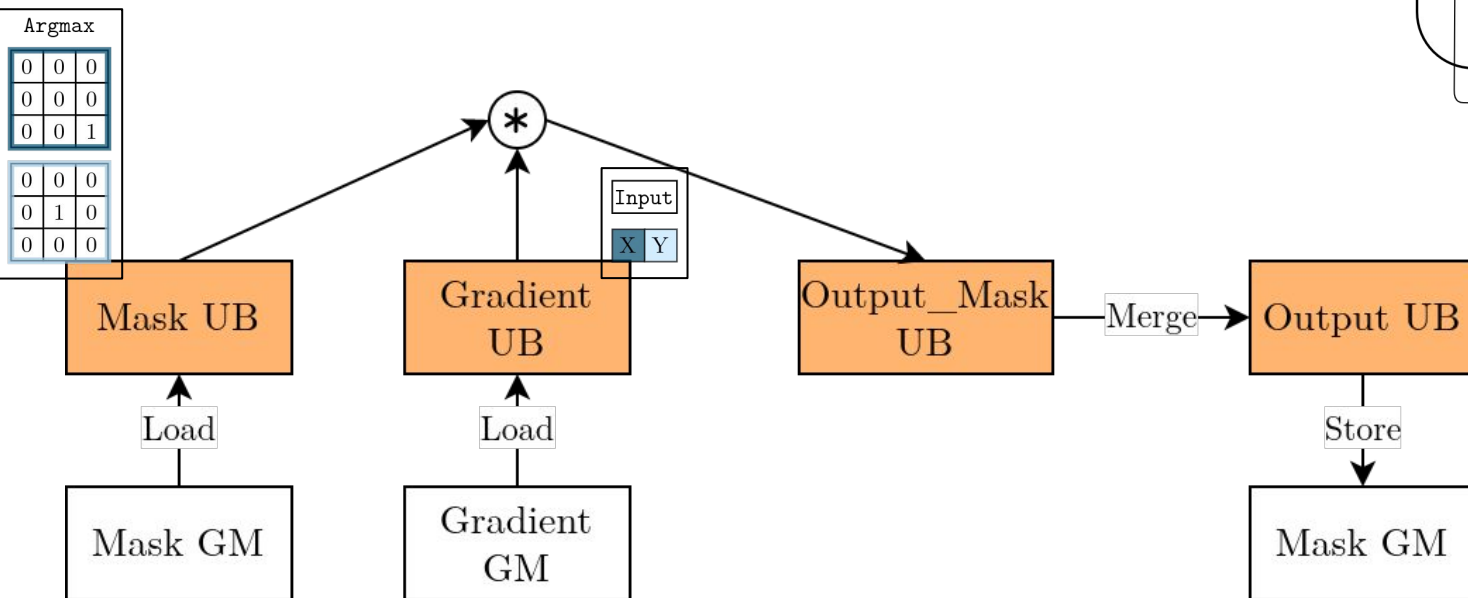
Im2col Based Pooling



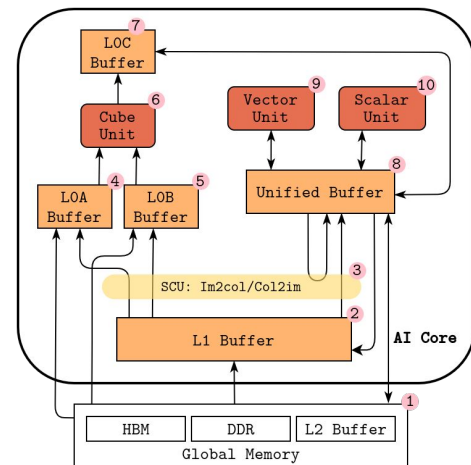
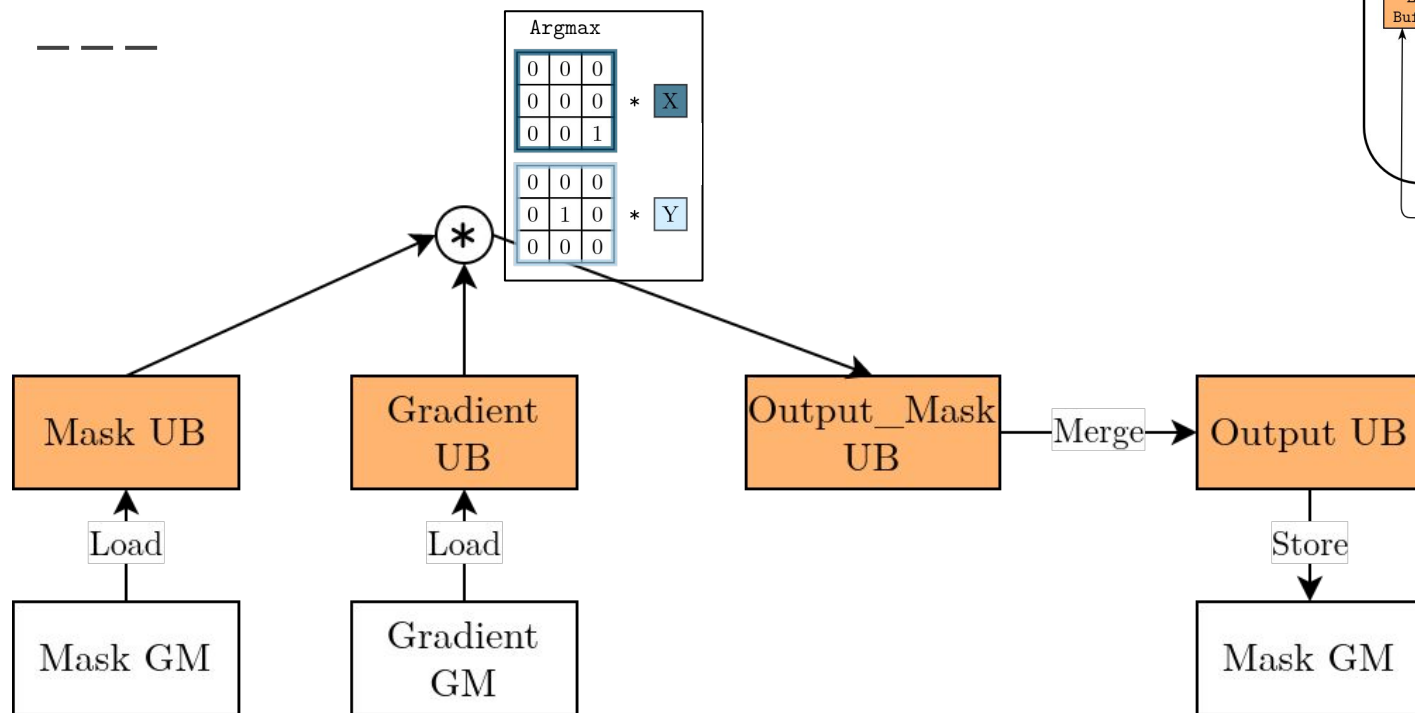
Backward Pooling for DaVinci



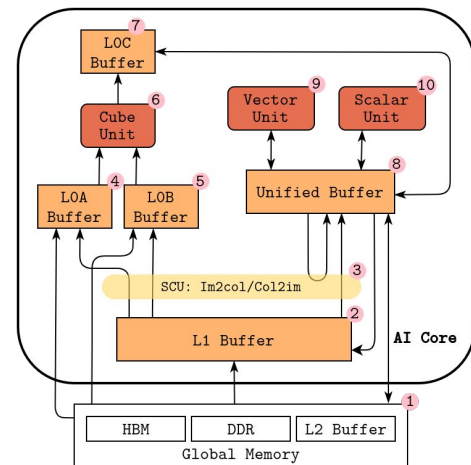
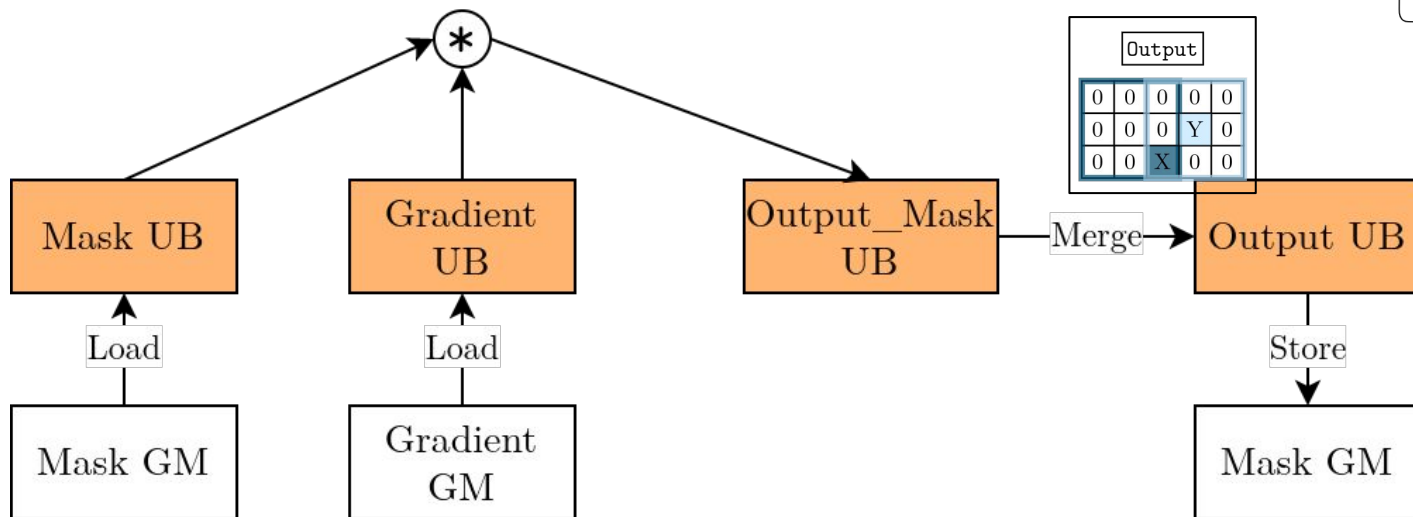
Backward Pooling for DaVinci



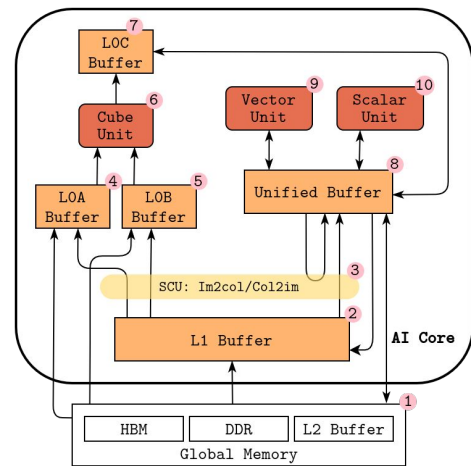
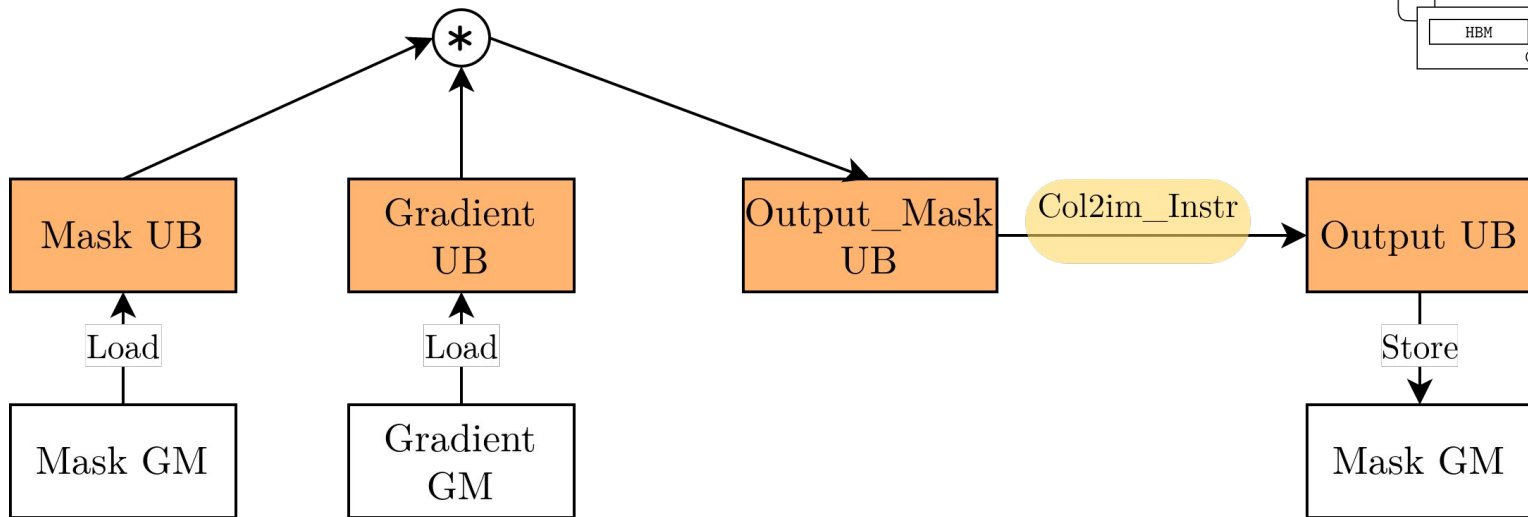
Backward Pooling for DaVinci



Backward Pooling for DaVinci

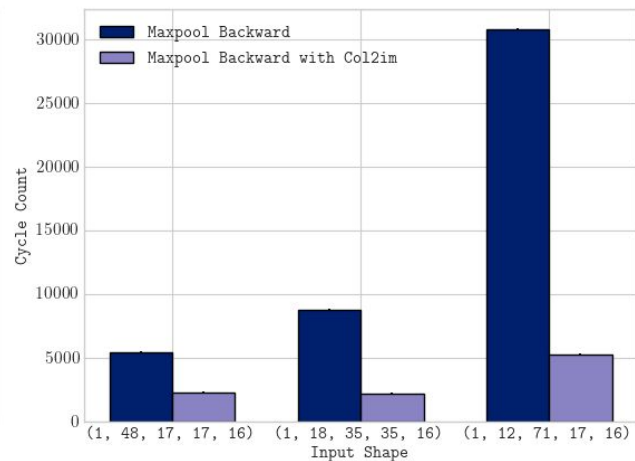
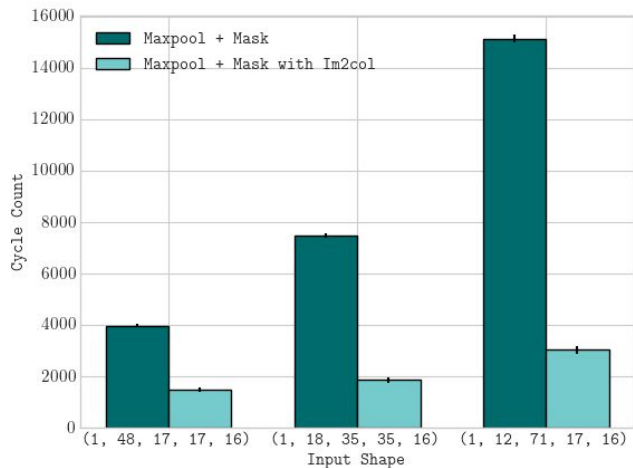
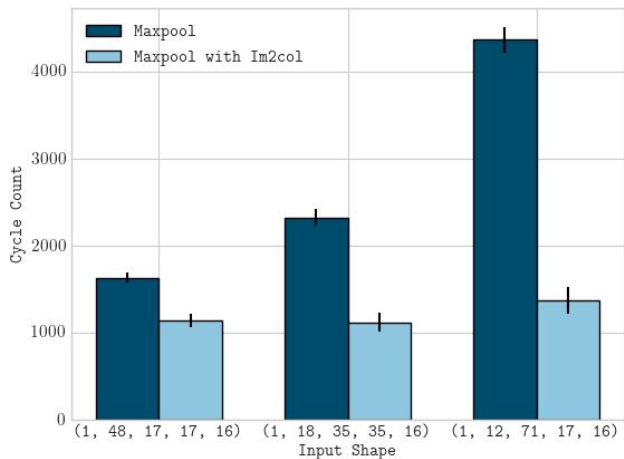


Col2im Based Backward Pooling



Maxpool Comparison

- Stride = (2, 2)
- Kernel size = (3,3)
- TVM

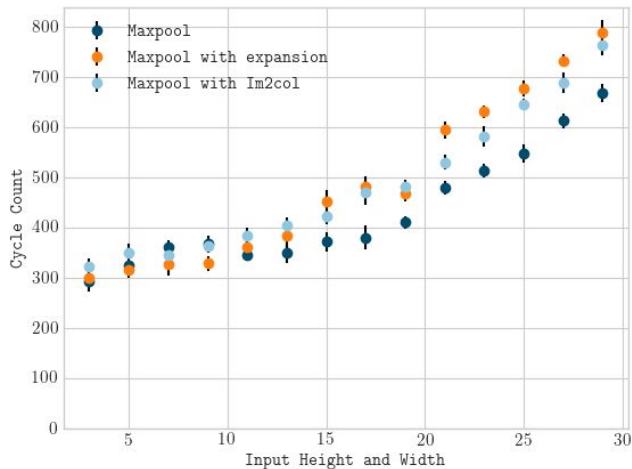


Maxpool Stride Comparison

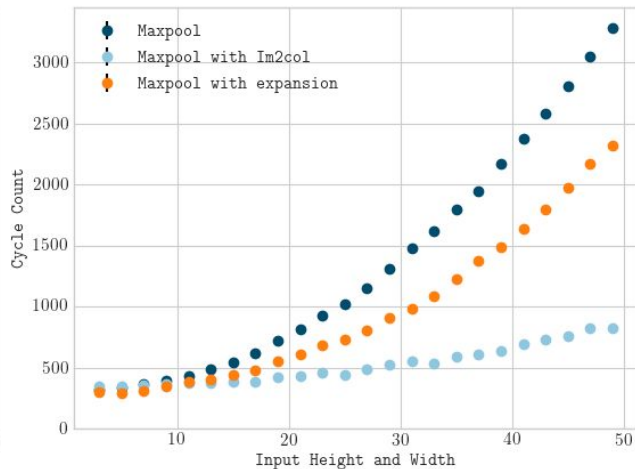
— — —

- Single channel (C_1)
- Kernel size = (3,3)
- Up to tiling threshold

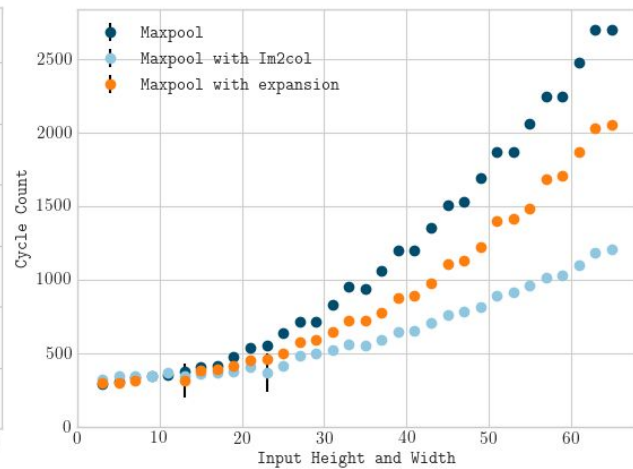
Stride = (1,1)



Stride = (2,2)



Stride = (3,3)



Conclusion

— — —

