

Pooling Acceleration in the DaVinci Architecture Using Im2col and Col2im Instructions

Caio S. Rohwedder⁺, João P. L. de Carvalho⁺, José Nelson Amaral^{*},
Guido Araújo⁺, Giancarlo Colmenares[†], Kai-Ting Amy Wang[†]

⁺*Institute of Computing
University of Campinas (UNICAMP)
Campinas, SP, Brazil
{joao.carvalho,guido}@ic.unicamp.br
{c157754}@dac.unicamp.br*

^{*}*Dept. of Computing Science
University of Alberta
Edmonton, AB, Canada
{jamaral}@ualberta.ca*

[†]*Heterogeneous Compiler Lab.
Huawei Canada Research Centre
Markham, ON, Canada
{kai.ting.wang}@huawei.com
{giancarlo.colmenares}@huawei.com*

Abstract—Image-to-column (Im2col) and column-to-image (Col2im) are data transformations extensively used to map convolution to matrix multiplication. These transformations rearrange the inputs of convolution to avoid its strided memory access pattern, thus providing a friendlier data layout for CPUs and GPUs. In artificial intelligence (AI) accelerators, these transformations allow convolution to be computed in matrix-multiplier units. Implemented in software, however, they impose a significant overhead that must be compensated by the efficiency gains of matrix multipliers. DaVinci is an AI accelerator architecture that introduces instructions to optimize Im2col and Col2im. Another core layer of convolutional neural networks that presents a strided memory access pattern is pooling. This paper explores the specialized Im2col and Col2im instructions to accelerate pooling layers in DaVinci. An experimental evaluation reveals that the proposed pooling implementations can yield speedups of up to 5.8 times compared to a baseline that does not use these specialized instructions. The speedups follow from an improved memory layout in the inputs of pooling, as this layout leads to better utilization of the vector processing unit in DaVinci.

Index Terms—CNN, AI Accelerator, Maxpool, Gradient, TVM

I. INTRODUCTION

With the increasing adoption of convolutional neural networks (CNNs), the optimization of each of its components has become fundamental. Convolution has been the main target of optimization because it is the most used and expensive layer in CNNs. However, many modern CNN architectures also use pooling to extract translation-invariant features and to perform subsampling. Max-pooling is the main variant of pooling that subsamples using the maximum value. While the performance impact of pooling is low compared to convolution, a naive implementation can hinder the overall performance of a CNN [1].

DaVinci [2] is an AI accelerator architecture that implements scalar, vector, and matrix multiplier units. The matrix multiplier unit allows efficient computation of convolution and other CNN layers, such as the fully connected, that can be mapped to matrix multiplication [3]. Convolution is mapped to matrix multiplication through the Im2col and Col2im data transformations. These transformations are memory-intensive and add significant performance overhead to convolution. However, highly optimized solutions for matrix multiplication

both in software (*e.g.*, OpenBLAS and Eigen libraries) and in hardware (*e.g.*, matrix multipliers) overcome this overhead. Still, DaVinci introduced instructions to optimize Im2col and Col2im. First, Im2col is performed during a load instruction (Im2Col) just before data reaches the memory buffers closest to DaVinci’s computational units. As such, this operation uses no temporaries and its memory overhead is only seen in these buffers. Second, Col2Im is vector instruction capable of better vectorizing over the scattered access pattern of Col2im. By using these instructions, convolution is computed in the matrix multiplier unit at a low overhead.

Max-pooling also has a strided access pattern, but unlike convolution it cannot be mapped to the matrix multiplier. Even so, its implementation can leverage the specialized Im2Col and Col2Im instructions. This paper thus proposes two key ideas to accelerate pooling in DaVinci: to produce an improved data layout by applying Im2Col instructions to the input of forward pooling, and to apply Col2Im instructions to the backward pooling instead of traditional vector instructions. Previous attempts to accelerate CNNs using FPGAs proposed pooling-specific instructions and computational units [4], [5]. Whereas the proposed approach uses a general-purpose vector computational unit and instructions primarily designed for convolution. Earlier work on improving pooling also overlooks its backward implementation [4], [6], [7], which is essential for training. Lastly, operation fusion, which effectively improves pooling paired with convolution [6], [8], is independent of the Im2col/Col2im based implementation presented in this work. Both optimizations can be applied in conjunction.

The main contributions of this paper are:

- A description of DaVinci’s Im2Col and Col2Im instructions, showing how they are executed and how they integrate into DaVinci’s datapaths (Section III).
- An approach to accelerate pooling with an Im2col-based forward implementation and a Col2im-based backward implementation using the DaVinci-specific Im2Col and Col2Im instructions (Section V).
- A rigorous evaluation of multiple pooling implementations in DaVinci, revealing speedups of up to 5.8 times on the Im2col/Col2im based implementations (Section VI).

83 The remaining sections are organized as follows: founda-
 84 tional concepts for this work appear in Section II. Section IV
 85 presents the software stack used to implement pooling opera-
 86 tors for DaVinci. Finally, Section VII presents related works,
 87 and Section VIII concludes this paper.

88 II. BACKGROUND

89 A. Convolution, Im2col and Matrix Multiplication

90 Convolution is a filtering operation used in image process-
 91 ing. This operation is the main building block of CNNs [9].
 92 In them, convolution repeatedly applies a kernel — a multi-
 93 channel filter composed of trainable weights — over patches of
 94 the input image. Patches are regions of the input that have the
 95 same size as the kernel. They are selected based on the stride
 96 parameters S_h and S_w , and given these parameters may or may
 97 not overlap. In an application of a kernel, its weights multiply a
 98 patch of the input. The multiplied results are summed together
 99 to generate a single output. A kernel is applied over each
 100 patch to generate a two-dimensional output, which is called
 101 a feature map. Convolution uses multiple kernels to produce
 102 multiple feature maps that are stacked as channels into a three-
 103 dimensional output.

104 The memory layout for the input of a convolutional layer is
 105 commonly described as $NCHW$, where each character repre-
 106 sents a dimension of a four-dimensional input: the number of
 107 images stacked together (N), channels (C), height (H), and
 108 width (W). The character's order specifies the order in which
 109 each dimension is arranged in memory. For simplicity, the
 110 dimension N has a length equal to one throughout the paper.

111 Convolution unrolling, also known as Image-to-Column
 112 (Im2col¹), is a data transformation that allows the mapping
 113 of convolution into matrix-matrix multiplication [11]. This
 114 transformation, illustrated in Figure 1, consists of creating two
 115 matrices, Out_{In} and Out_{Ker} , based on the input image and
 116 the kernels, respectively. Each row of matrix Out_{In} contains
 117 all the input needed to compute one element of an output
 118 feature map linearized into one dimension. Each column of
 119 matrix Out_{Ker} contains the weights of a kernel similarly
 120 linearized. Thus, multiplying Out_{In} and Out_{Ker} is equivalent
 121 to performing convolution with its original inputs.

122 If the stride sizes (S_h, S_w) are smaller than the kernel's
 123 height and width (K_h, K_w), patches will overlap. The over-
 124 lapping elements will be copied to multiple rows of matrix
 125 Out_{In} , resulting in a bigger memory footprint. This is the
 126 main drawback of the Im2col technique when contrasted with
 127 direct-convolution based approaches. An example with a single
 128 channel is shown in Figure 2. The two patches are highlighted
 129 and they overlap on the elements $\{3, 8, 13\}$. As a result, these
 130 elements appear in both rows of the output of Im2col (on
 131 the right). Nonetheless, Im2col is used across AI frameworks
 132 to implement convolution because matrix multiplication offers
 133 an input with a friendlier memory layout to CPUs and GPUs,

¹The transformation of the input image can also be an image-to-row transformation if the multiplication is transposed $(AB)^T = B^T A^T$ [10]. The Im2col name will be used to refer to all variants of this transformation.

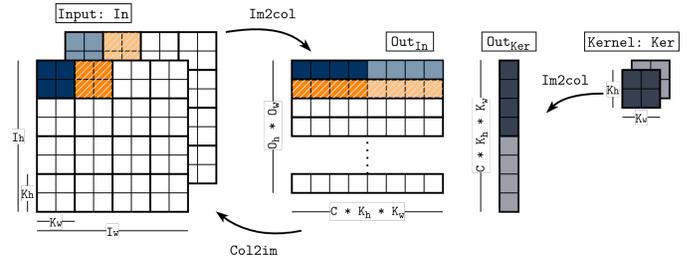


Fig. 1. Im2col: $In (C, I_h, I_w)$ is transformed into matrix $Out_{In} (O_h \times O_w, C \times K_h \times K_w)$ and a single kernel (C, K_h, K_w) is transformed into the matrix $Out_{Ker} (C \times K_h \times K_w, 1)$, where O_h and O_w represent the number of patches in the height and width of the input. The bold squares (2, 2) in In represent patches of the image to which the kernel is applied. Col2im: the backward operator of Im2col, from Out_{In} to In .

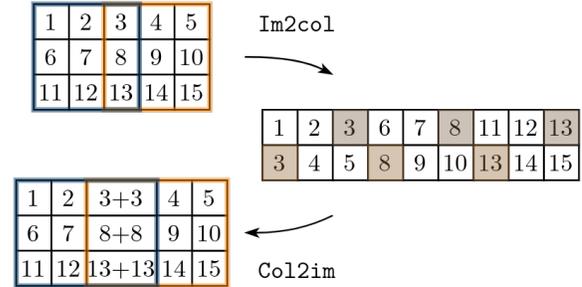


Fig. 2. Im2col and Col2im performed on two overlapping patches.

134 making it easier to apply vectorization techniques [11]. The
 135 availability of optimized linear algebra libraries such as Open-
 136 BLAS [12], ATLAS [13], and Eigen [14], and AI accelerator
 137 designed around matrix multiplier units, further incentivizes
 138 such a transformation.

139 B. Backward Operators and Col2im

140 To train a neural network, the input values are first propa-
 141 gated forward to produce an output. Then, the error between
 142 this generated output and the expected output is calculated
 143 through a loss function. The gradient of this loss function is
 144 propagated backward towards the input so that the network
 145 can be tuned. Thus, every forward operator has a dual-
 146 operator applied in the backward pass, namely its *backward*
 147 *operator* [15].

148 The backward operator of Im2col is called Col2im, and it
 149 is also illustrated in Figure 1. Col2im is used in the backward
 150 propagation pass of convolutional layers implemented with
 151 Im2col. The incoming gradients in the shape of the matrix
 152 Out_{In} are propagated back to the original $NCHW$ layout. If
 153 there is no overlap, as in the example of Figure 1, Col2im
 154 simply returns the matrix to its original shape. But when
 155 patches do overlap, gradients that refer to the same position
 156 in the output are summed, as shown in Figure 2.

157 C. Pooling Operators

158 Spatial feature pooling subsamples images to obtain
 159 translation-invariant feature maps in computer-vision archi-

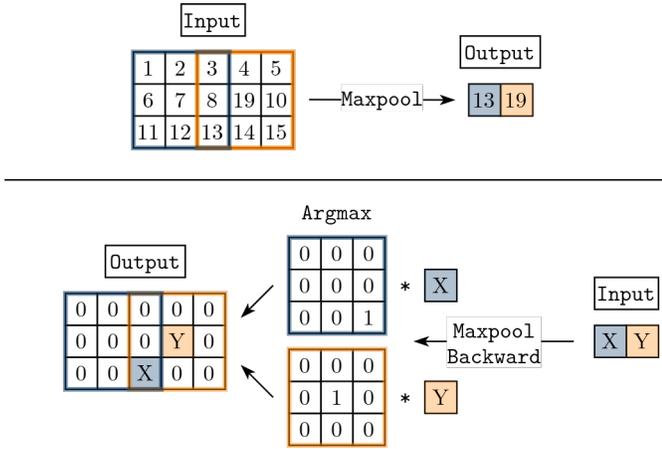


Fig. 3. Forward and backward computation of MaxPool for two overlapping patches.

160 tures [16]. Similar to convolution, pooling is one of the
 161 building blocks of CNNs. Pooling layers are commonly used in
 162 modern CNN architectures such as Resnet [17], Inception [18],
 163 and Xception [19]. Pooling also applies a kernel over patches
 164 of its $NCHW$ input. But unlike convolution, the kernel has no
 165 weights, it only selects patches based on the stride parameters.
 166 A reduction function is applied to the selected patches to
 167 subsample the input. This reduction is typically applied to
 168 the height (H) and width (W) dimensions of the input,
 169 operating on the channels independently. As a result, the output
 170 of pooling has the same number of channels as the input.
 171 Different reduction functions can be chosen: the max function
 172 selects the maximum value (MaxPool), and the avg function
 173 computes the average of the patch (AvgPool). MaxPool is
 174 preferred among CNNs as it looks at the maximal activation
 175 of features, rather than diluting them with an average [20].
 176 Figure 3 (on top) shows an example of MaxPool forward.

177 The implementation of backward pooling depends on the
 178 reduction function utilized. For Maxpool, each input is mul-
 179 tiplied by its corresponding Argmax mask, where the position
 180 of the maximum element in the original patch is set to 1 and
 181 all the other positions are set to 0. Next, as with Col2im,
 182 the masks return to the original $NCHW$ shape and the over-
 183 lapping elements are summed together. This output correlates
 184 how much a change in each input element of MaxPool forward
 185 affects its output elements [21]. Figure 3 shows an example of
 186 Maxpool backward on two overlapping patches. In summary,
 187 the gradients are only propagated backward to the maximum
 188 elements [22].

189 III. THE DAVINCI ARCHITECTURE

190 DaVinci [2] is an AI accelerator architecture used by
 191 Huawei’s Ascend chips. The following subsections describe
 192 components of the Ascend 910 chip.

193 A. AI Core

194 Figure 4 shows a closer view of DaVinci’s main component,
 195 the AI Core, and its corresponding data paths. The AI Core is

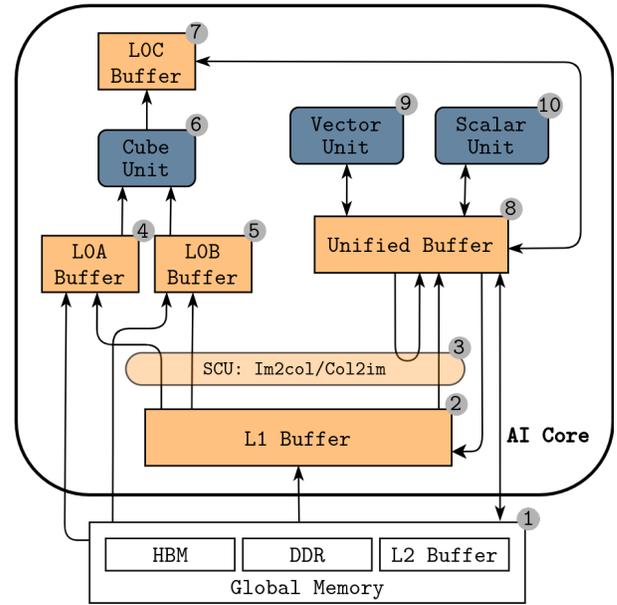


Fig. 4. Data paths of the AI Core.

196 composed of three processing units (Cube, Scalar, and Vector
 197 Unit), a set of private buffers (LOA, LOB, LOC, L1, and Unified
 198 Buffer), and a Storage Conversion Unit (SCU). Outside of the
 199 AI Core sits the Double Data Rate (DDR) and High Bandwidth
 200 Memory (HBM) memories and an L2 Buffer, all of which are
 201 shared among the AI Cores of a chip.

202 Both Scalar and Vector Units operate on data loaded from/
 203 stored to the Unified Buffer. The Vector Unit performs basic
 204 arithmetic and logic vector operations (e.g., subtracting two
 205 vectors). It uses a 128-bit mask register in which every bit
 206 represents one element of a vector instruction that may be
 207 processed or not. The Scalar Unit has both general and special-
 208 purpose registers, which are used to execute control-flow and
 209 scalar arithmetic operations, as well as index and address
 210 calculations.

211 The Cube Unit is based on a multidimensional systolic
 212 array [23], it implements matrix multiplication using an array
 213 of processing elements that perform multiply-accumulate op-
 214 erations. This unit acts similarly to the Matrix Multiplier Unit
 215 (MXU) of Google’s Tensor Processing Unit [24]. Buffers LOA
 216 and LOB store the inputs of the Cube Unit, and the LOC buffer
 217 stores its output. While the operands for the Vector Unit are
 218 vectors, the Cube Unit receives *data-fractals* from its input
 219 buffers. A *data-fractal* is a small matrix of a constant size of
 220 4096 bits. The Cube Unit can multiply two data-fractals per
 221 clock cycle.

222 The private buffers of the AI Core (LOA, LOB, LOC, L1, and
 223 Unified Buffer) are organized as scratch-pad memories [25].
 224 Data movement between these buffers must be explicitly man-
 225 aged by the application, in contrast, hardware-managed caches
 226 are transparent to the application and ensure consistency by
 227 hardware protocols. Thus, the programmer needs to specify
 228 which data should be brought to each buffer, and also needs

to maintain data consistency. In a scratch-pad memory, each buffer has its own address space, which is separated from the address space of the memory. With this organization, more complexity is placed upon the application’s code, but it comes with the benefit of not requiring tag bits, dirty bit, and the comparison logic that transparent caches need in the hardware. From the AI Core’s perspective, all shared memories (DDR, HBM, and L2) are considered global memory and are represented as ① in Figure 4. Given that their data-paths are the same, they are drawn only once.

The Storage Conversion Unit (SCU) may perform many data-layout transformations when data is transferred between buffers. This unit implements Im2col, Col2im, and other transformations, out of the scope of this work, such as padding, matrix-tile transposition, and sparse-matrix decompression. The SCU enables instructions, such as Im2Co1, to perform fast layout transformations while data is transferred between buffers. As a result, the memory overhead that these transformations may imply appears only on the target buffers. Such instructions were specifically designed to operate on the memory layout described next.

B. Fractal Memory Layout

To avoid memory alignment and padding problems in the Cube Unit, DaVinci includes the constant-length dimension C_0 in the representation of an input image. As a result, a slight variation of $NCHW$ is used, called the *fractal memory layout*. This format is represented by NC_1HWC_0 , in which C_0 represents part of a split in the channel dimension (C) of $NCHW$. To make the conversion from $NCHW$ to NC_1HWC_0 , C is split into C_1 and C_0 , where $C_1 = \lceil C/C_0 \rceil$. If the original number of channels (C) is not divisible by C_0 , the C_0 dimension must be zero-padded to reach its required length. Given a data type, the length of C_0 makes the inputs of the Cube Unit (data-fractals) always have 4096 bits of data. A data-fractal has $16 * C_0$ elements, thus for `Float16`, C_0 has a length of 16. For `Unsigned8`, C_0 has a length of 32. The data type `Float16` is adopted in this paper.

C. Im2Col Instruction

Im2Co1 is a data-transformation instruction executed in the SCU that acts as a load instruction. It may be applied to a data-fractal that is loaded from L1 to LOA ②→④ and LOB ②→⑤, so as to prepare data for computation in the Cube Unit. It may also be applied to a data-fractal that is loaded from L1 to the Unified buffer ②→⑧, to prepare data for computation in the Vector and Scalar Units.

There are two main differences when comparing the Im2Co1 instruction to the Im2col transformation shown in Figure 1. First, Im2Co1 is a single instruction, it is only able to load and transform one fractal of an image at a time. Even if it could operate on a whole image, its target buffers (LOA, LOB, Unified Buffer) may not be capable of storing the transformed image. For this reason, Im2Co1 instructions can be used to load and transform a tile of an input. Second, Im2Co1 is designed to load an input that is in the fractal memory

layout NC_1HWC_0 . Therefore, its output will also have a different memory layout when compared to the one shown on the right of Figure 1. The advantage of performing Im2col as load instruction is that the increase in memory overhead from duplicated elements only appears in the target buffers (LOA, LOB, and Unified buffer), which are the buffers closest to the Cube and Vector Units.

Im2Co1 needs the following parameters related to the input image (or tile), which are constant for all instructions loading the same input:

- Height (I_h) and width (I_w) of the input image;
- Left (P_l), right (P_r), top (P_t), and bottom (P_b) zero padding;
- Stride in the height (S_h) and width (S_w) directions;
- Kernel height (K_h) and width (K_w).

Based on these parameters, the number of patches (O_h, O_w) in the input’s height and width can be calculated by Equation 1. Furthermore, each Im2Co1 instruction needs the three following positional parameters to choose which elements of the input it will load, in which the parameters (x, y) are coordinates in the height and width (HW) dimensions of the input.

- The starting position in the image (x, y);
- Relative position inside of a patch (x_k, y_k);
- Access index of the C_1 dimension (c_1).

$$\begin{aligned} O_h &= \left\lfloor \frac{I_h + P_b + P_t - K_h}{S_h} \right\rfloor + 1 \\ O_w &= \left\lfloor \frac{I_w + P_l + P_r - K_w}{S_w} \right\rfloor + 1 \end{aligned} \quad (1)$$

To load a fractal (16 rows of C_0 elements) to a buffer, Im2Co1 performs the following tasks: (i) process each element of dimension N individually; (ii) access the element c_1 of dimension C_1 ; (iii) select the next 16 consecutive patches starting from position (x, y); (iv) select the elements in the (x_k, y_k) position, relative to each of the 16 patches; (v) load the C_0 dimension for the 16 selected elements; (vi) store the loaded elements as a fractal into the target buffer.

Figure 5 exemplifies a small image loaded using four Im2Co1 loads. The input image is in the fractal layout NC_1HWC_0 , but the lengths of N and C_1 are 1, so they are not shown. The parameters used in this example correspond to: $(I_h, I_w) = (8, 8)$, $(K_h, K_w) = (2, 2)$, $(S_h, S_w) = (2, 2)$, and $(O_h, O_w) = (4, 4)$. Notice that there is no padding. The input has exactly 16 patches (bold squares), so (x, y) is set to the first position (0, 0) and is not changed afterward. For the first Im2Co1 (blue squares), $(x_k, y_k) = (0, 0)$, while for the second (orange squares), $(x_k, y_k) = (0, 1)$. Two more Im2Co1 instructions are issued, corresponding to (x_k, y_k) equal to (1, 0) and (1, 1). This results in four fractals concatenated side by side. If there were more patches in the image, (x, y) would be changed to another position to create a new row of fractals in the output. Bigger inputs are loaded by issuing multiple Im2Co1 instructions while iterating the positional parameters sequentially. This iteration can be seen as if it composed a

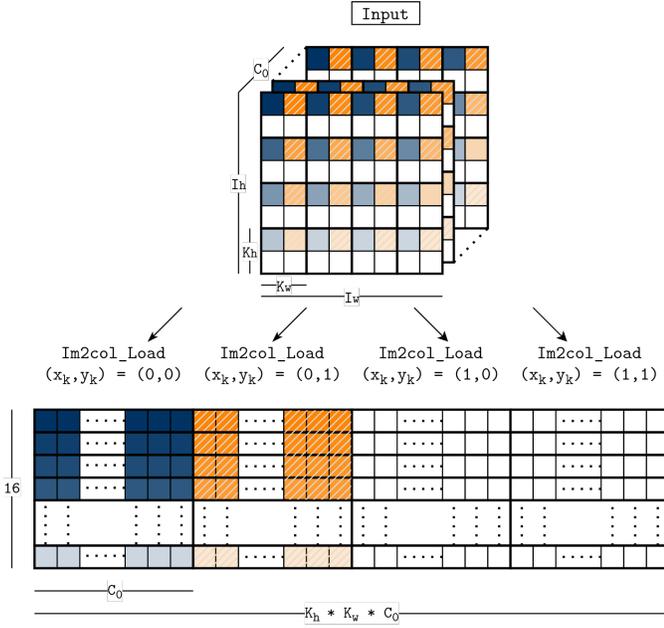


Fig. 5. Four Im2Co1 loads. The input is HWC_0 . On the bottom, are the four resulting fractals of size $16 \times C_0$. The difference between the loads is the position relative to the patch (x_k, y_k) , which is $(0, 0)$ for the first load (highlighted in blue) and $(0, 1)$ for the second (highlighted in orange), $(1, 0)$ for the third, and $(1, 1)$ for the fourth. The resulting fractals are concatenated in the output buffer.

333 triple-nested loop with iterator vector in the form of $[(x, y),$
 334 $c_1, (x_k, y_k)]$, from the outermost to the innermost loop.

335 As with most instructions in the DaVinci architecture,
 336 Im2Co1 supports a repetition parameter that causes an in-
 337 struction to be reissued automatically. For Im2Co1 there are
 338 two possible repetition modes. Mode 0 repeats Im2Co1 for
 339 the next positions inside the kernel (x_k, y_k) , from $(0, 0)$ to
 340 $(0, 1)$, for example. If the length of C_1 is bigger than 1,
 341 Im2Co1 in repetition mode 0 will continue to the next c_1
 342 index and iterate over (x_k, y_k) again. This repetition mode
 343 acts as the loops of $[c_1, (x_k, y_k)]$, but multiple Im2Co1 are
 344 needed to also change (x, y) . Therefore, the input in Figure 5
 345 can be fully loaded by issuing a single Im2Co1 starting at
 346 $(x_k, y_k) = (0, 0)$ with repeat mode 0 to repeat four times,
 347 changing (x_k, y_k) from $(0, 0)$ to $(0, 1)$, $(1, 0)$ and $(1, 1)$. Mode
 348 1 reissues Im2Co1 for the next (x, y) position after skipping
 349 the 16 currently selected patches. In this mode, one Im2Co1
 350 instruction acts as the loop of $[(x, y)]$, and multiple instructions
 351 are needed to change c_1 and (x_k, y_k) , thus, (x, y) becomes
 352 the innermost loop of the iterator vector. If the nesting order
 353 of these loops changes, so does the order in which fractals
 354 are stored in memory. By changing the order from $[(x, y),$
 355 $c_1, (x_k, y_k)]$ to $[c_1, (x_k, y_k), (x, y)]$ in mode 1, Im2Co1 will
 356 store fractals in a transposed order resulting in an output
 357 matrix of shape $(C_1 \times K_h \times K_w \times 16, (O_h \times O_w)/16 \times C_0)$.
 358 This shape can also be considered as a tensor of dimensions
 359 $(C_1, K_h, K_w, O_h, O_w, C_0)$, which is the shape used in the
 360 accelerated forward pooling implementation in Section V.

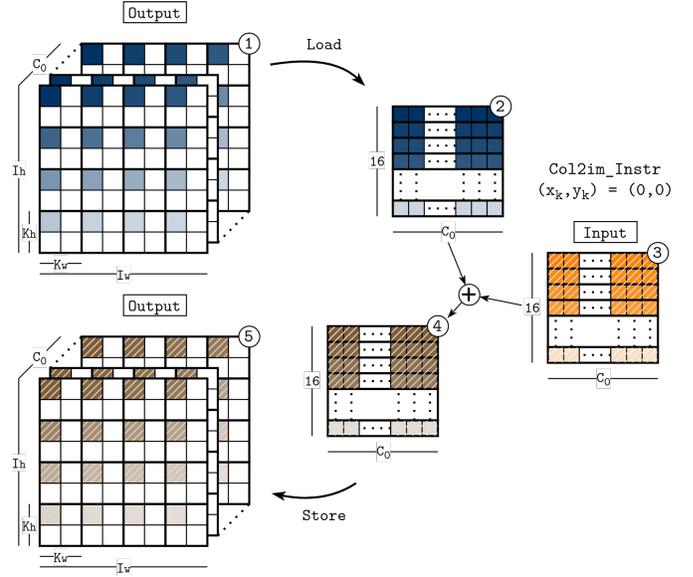


Fig. 6. Single Col2Im load with parameters $(x, y) = (0, 0)$ and $(x_k, y_k) = (0, 0)$.

D. Col2Im Instruction

361 Col2Im is an instruction that is used as the backward
 362 operator of Im2Co1. It acts as a vector instruction that loads
 363 data from and stores data to the Unified Buffer $\textcircled{8} \rightarrow \textcircled{8}$
 364 (Figure 4). Col2Im takes fractals as inputs and stores them
 365 in the NC_1HWC_0 format. Because of this, Col2Im receives
 366 the same parameters as Im2Co1 referring to its output. Besides
 367 the change in memory layout, if two patches overlap in the
 368 output, input elements that refer to the same output position
 369 need to be summed. This sum is shown in Figure 2, but it is
 370 performed at an instruction level. For that, Col2Im requires
 371 its output to be initialized with zeros.
 372

373 Figure 6 shows how a single Col2Im instruction works
 374 with an already initialized output. This example uses the same
 375 parameters as the first (blue) Im2Co1 shown in Figure 5. In
 376 Figure 6, Col2Im loads the initialized output $\textcircled{1}$ in an Im2Co1
 377 manner $\textcircled{2}$. Then, it sums the loaded fractal with the input
 378 fractal $\textcircled{3}$. Finally, it stores the resulting fractal $\textcircled{4}$
 379 back to its corresponding positions in the output $\textcircled{5}$. This example could
 380 not be loaded using a repetition because the only repetition
 381 mode available for Col2Im is mode 1. It works as in Im2Co1
 382 by changing the (x, y) parameters and thus requires an input
 383 with more than 16 patches.

IV. SOFTWARE STACK

384 A C-like language called CCE (Cube-based Compute En-
 385 gine) C is used to write code for DaVinci chips. Because it
 386 is a very low-level language, implementing and optimizing
 387 multiple AI operators manually is a cumbersome and error-
 388 prone task. The Automatic Kernel Generator (AKG), a tool
 389 for operator design and also a library of operators, is used
 390 to enable the design of AI operators in CCE C. AKG uses
 391 TVM's [26] domain-specific language (DSL) to design its
 392

393 operators, which are lowered to CCE C by its compiler passes.
 394 For every operator that is defined with AKG, its backward
 395 operator is also needed to allow training.

396 A. Scheduling for DaVinci

397 TVM’s DSL is based on the Halide language [27]. The main
 398 idea of both languages is to decouple the execution definition
 399 (the algorithm) from the execution strategy (the algorithm’s
 400 schedule). With this separation, the programmer is free to
 401 test multiple optimization strategies by rewriting a schedule
 402 without changing the algorithm. The schedule allows the use of
 403 techniques such as function inlining and loop transformations
 404 (e.g., tiling, fusion, unrolling, and loop vectorization). The
 405 decoupling of the algorithm from its schedule is possible
 406 because Halide’s and TVM’s DSLs are tailored respectively
 407 for image processing and deep learning algorithms. There
 408 is a high degree of data parallelism in applications from
 409 these fields [27] as their algorithms are mainly composed
 410 of loops with no dependencies between iterations, known as
 411 DOALL loops [28]. In this scenario, the loop transformations
 412 previously mentioned are trivial.

413 TVM allows code generation for other backends besides
 414 CPUs. Hence, schedules can explicitly refer to a backend-
 415 specific construct. For example, schedules allow binding loops
 416 in the algorithm to blocks and threads, which are constructs
 417 found in GPUs. AKG uses the same principle to generate
 418 code for DaVinci devices. A DaVinci-specific schedule is
 419 responsible for controlling the movements of data between
 420 the scratch-pad buffers and for specifying computations that
 421 are local to a buffer. Together with the backend-specific
 422 schedule primitives, it is possible to apply other optimization
 423 techniques (e.g., tiling) to improve the locality of memory
 424 accesses. Between all the possible primitives, two are handled
 425 automatically by AKG: vectorization and parallelization. First,
 426 the inner loops of computations are vectorized (minimally on
 427 the C_0 dimension) so that the Vector Unit is utilized automat-
 428 ically. When possible, the vector instructions are also issued
 429 with repeat factors. Second, the outer loops are parallelized
 430 between the AI Cores available on the target device. These
 431 default behaviors are similar to those taken by Halide’s auto-
 432 scheduler [29]. AKG also has a polyhedral framework that
 433 automatically schedules computation on DaVinci, but it does
 434 not support all instructions (e.g., Col2Im).

435 V. IM2COL/COL2IM BASED POOLING

436 The Vector Unit computes pooling in DaVinci. The perfor-
 437 mance of vector instructions running in it depends mostly on
 438 two factors. First, the vector mask should be saturated so that
 439 all vector lanes are utilized and parallelism is maximized. Sec-
 440 ond, the repetition parameter should be employed, thus remov-
 441 ing loops and barriers around vector instructions, and taking
 442 pressure off instruction fetching. Ideally, a single instruction
 443 should operate over an entire tensor (or tile) present in the
 444 Unified buffer. This Section describes the Im2col/Col2im
 445 based pooling implementations in comparison to their standard

```

1 input = placeholder((N, C1, Ih, Iw, Co),
2                     name="input")
3 red_h = reduce_axis((0, Kh), "red_h")
4 red_w = reduce_axis((0, Kw), "red_w")
5 output = compute((N, C1, Oh, Ow, Co),
6                 lambda n, c1, h, w, c0:
7                 max(input[n, c1,
8                       h*S_h+red_h,
9                       w*S_w+red_w,
10                      c0],
11                    axis=[red_h, red_w]))

```

Listing 1. MaxPool defined with TVM’s DSL

446 implementations in TVM. Lowered CCE C code is used to
 447 highlight the above-mentioned factors in each implementation.

448 A. Maxpool Forward

449 A standard TVM implementation of Maxpool forward is
 450 represented in Listing 1. It describes its input and output
 451 shapes (Lines 1 and 5, respectively) and its 2D max reduction
 452 of each patch (Lines 6 to 10). Using TVM’s schedule, this
 453 computation is divided in the C_1 dimension so that a tile
 454 of size (I_h, I_w, C_0) is computed at a time, producing an
 455 output tile of size (O_h, O_w, C_0) unless further tiling is needed.
 456 Each input tile is sent from global memory to the Unified
 457 Buffer of an AI Core ①→⑧. If multiple AI Cores are
 458 available, multiple tiles can be processed in parallel. After
 459 the computation is finished, the output tile is brought back to
 460 global memory ⑧→①.

461 This implementation is lowered to CCE C code where the
 462 vmx instruction is executed. vmx computes the maximum
 463 between elements of the output and input tiles and writes back
 464 to the output tile. For that, the output tile is initialized with
 465 the minimum value of the data type in use. In this setting, only
 466 16 of 128 elements of the vector mask are set, accounting for
 467 the innermost dimension C_0 of the tiles. Additionally, each
 468 vmx uses repetition to obtain the maximum value across the
 469 width of a patch K_w (the innermost reduction axis red_w). The
 470 vmx instruction is issued $O_h * O_w * K_h$ times to complete
 471 the computation. These suboptimal parameters result from the
 472 strided access pattern seen in Lines 8 and 9 of Listing 1.

473 The Im2col based implementation is described in Listing 2.
 474 Lines 1 and 3 represent the Im2col transformation. In this
 475 implementation, the input starts in the global memory with
 476 its original shape, which is tiled along the C_1 dimension.
 477 Next, the input is first loaded to the L1 buffer of an AI
 478 Core ①→② and then loaded with Im2Col to its Unified
 479 buffer using the repeat mode 1 ②→⑧. The transformed
 480 input has the shape shown in Line 3. Its tiles have a shape
 481 $(K_h, K_w, O_h, O_w, C_0)$ in the Unified buffer. The max re-
 482 duction now occurs in the outer (K_h, K_w) dimensions, as
 483 shown in Lines 10 and 11. Considering the input and output
 484 tiles, $(K_h, K_w, O_h, O_w, C_0)$ and (O_h, O_w, C_0) respectively,
 485 the lowered CCE C code is able to set all 128 elements
 486 of the vector mask, and, in conjunction with the repetition
 487 parameter, a single vmx computes the max between the entire
 488 output tile and the three innermost dimensions of the input tile,
 489 which are identical. This instruction is only issued $K_h * K_w$

```

1 input = placeholder((N, C1, Ih, Iw, C0),
2                   name="input")
3 input-ub = placeholder((N, C1, Kh, Kw, Oh, Ow, C0),
4                      name="input-ub")
5 red_h = reduce_axis((0, Kh), "red_h")
6 red_w = reduce_axis((0, Kw), "red_w")
7 output = compute((N, C1, Oh, Ow, C0),
8                lambda n, c1, h, w, c0:
9                  max(input-im2col[n, c1,
10                          red_h,
11                          red_w,
12                          h, w, c0],
13                    axis=[red_h, red_w]))

```

Listing 2. MaxPool performed on an input with the resulting shape of using the Im2Col load

```

1 argmax-mask = placeholder((N, C1, Kh, Kw, Oh, Ow, C0))
2 gradients = placeholder((N, C1, Oh, Ow, C0))
3 mask-gradient = compute((N, C1, Kh, Kw, Oh, Ow, C0),
4                        lambda n, c1, kh, kw, oh, ow, c0:
5                          argmax-mask(b, c1, kh, kw, oh, ow, c0)
6                          * gradient(b, c1, oh, ow, c0)
7                          )

```

Listing 3. Part of MaxPool backward defined with TVM’s DSL

490 times to finish the computation, effectively improving upon
491 the standard implementations of Listing 1.

492 For training, it is useful to save an additional result in the
493 forward implementation of Maxpool: the argmax mask. This
494 mask is used by Maxpool’s backward operator to store the
495 position of the maximum element of each patch, as shown
496 in Figure 3. This result is obtained by comparing each patch
497 of the input with its maximum value. Saving this mask is
498 independent of the use of Im2Col instructions. Still, the
499 Im2Col output shape of Line 3 in Listing 2 is used to store
500 it, as it keeps overlapping patches separated. This shape also
501 enables Maxpool backward to use Col2Im instructions, which
502 is described next.

503 B. Maxpool Backward

504 Maxpool backward receives two inputs: the argmax mask
505 and the incoming gradients. Listing 3 shows part of its
506 implementation. Line 3 defines a computation that multiplies
507 the patches in the argmax mask with their corresponding
508 gradients. This multiplication is represented on the bottom of
509 Figure 3. Next, the multiplied patches need to be merged back
510 into the original (N, C_1, I_h, I_w, C_0) shape by summing values
511 in the overlapping areas, which is not shown in Listing 3
512 for brevity. This merge step is critical for performance and
513 it is depicted on the bottom-left of Figure 3. Its TVM imple-
514 mentation requires expanding mask-gradient to a shape of
515 $(N, C_1, I_h, I_w, O_h, O_w, C_0)$, where each patch is copied only
516 once in its correct position in I_h and I_w , and other elements are
517 set to zero. The expanded representation is then reduced with
518 sum on dimensions O_h and O_w , effectively summing up the
519 overlapping areas in every patch and obtaining the final shape
520 of (N, C_1, I_h, I_w, C_0) . This expansion would be incredibly
521 costly due to its size, however, TVM allows it to be inlined
522 using a schedule. As a consequence, the patches are merged,
523 and the overlapping regions are summed directly to the final
524 output shape (from the shape $(N, C_1, K_h, K_w, O_h, O_w, C_0)$
525 to (N, C_1, I_h, I_w, C_0)). Besides the mentioned inlining, the
526 schedule works similarly to Maxpool forward, loading both
527 inputs from the global memory to the Unified buffer ①→⑧
528 so that the multiplication is computed in the Vector Unit ⑨,
529 and tiling the computation on C_1 .

530 The lowered code uses vmul for the multiplication step,
531 and vadd, for the merge step. These instructions work in the

same way as vmax, but for multiplication and addition. While
vmul works well in multiplying tiles of the gradient with the
mask, the scattered access pattern of the merge step leads
to very poor usage of the Vector Unit. That is because the
vadd instructions only set 16 elements of the vector mask
(vectorizing on C_0) and repetition is not used.

The Col2im based implementation comes from the obser-
vation that the merge step computes exactly the Col2im
operation. As mentioned before, Col2Im uses the Unified
buffer to load its input and to store its output ⑧→⑧.
Therefore, it is possible to use Col2Im instead of vadd. The
Col2Im instruction is able to load and store to the scattered
elements of the output, summing two fractals at a time, as
shown in Figure 6. In comparison with vadd that had 16 (C_0)
elements of the vector mask set, Col2Im enables vectorization
over $16 * 16$ elements (a fractal) at a time, and its repetition
mode can be used to operate over the entire tile in the Unified
buffer. A Col2Im instruction needs to be issued $K_h * K_w$ times
to complete the merge step of a tile. Therefore, switching vadd
for Col2Im presents a good opportunity for performance gains.

552 C. Avgpool

553 The forward and backward operators of Avgpool are similar
554 to those described before. But opposed to Maxpool, the
555 forward implementation reduces using sum instead of max.
556 Consequently, its CCE C code uses vadd instead of vmax. But
557 regardless of this change, the access pattern stays the same and
558 can benefit from using Im2Col. Additionally, a new operation
559 is needed to compute an element-wise division before saving
560 the final output. As for the backward operator, there is no need
561 to use the Argmax mask as an input. The equivalent mask for
562 Avgpool contains 1 in all its positions, given that all input
563 elements contribute to the output of a sum. Besides the mask,
564 the backward implementation is the same and it can also use
565 Col2Im instructions.

566 VI. EXPERIMENTAL EVALUATION

567 This evaluation compares the performance of the
568 Im2col/Col2im based Maxpool with the standard TVM
569 Maxpool implementation described in Section V. All the
570 experiments were run on an Ascend 910 chip, which contains
571 32 AI Cores. The cycle count numbers were obtained using
572 the hardware counters of the chip, and they refer to the on-
573 chip execution time running at a frequency of 100 MHz. The
574 cycle count is currently the only metric that could be obtained
575 from the chip. Each evaluation was repeated ten times, and
576 the graphs show the average value and a 95% confidence
577 interval. To use the Im2Col and Col2Im instructions in

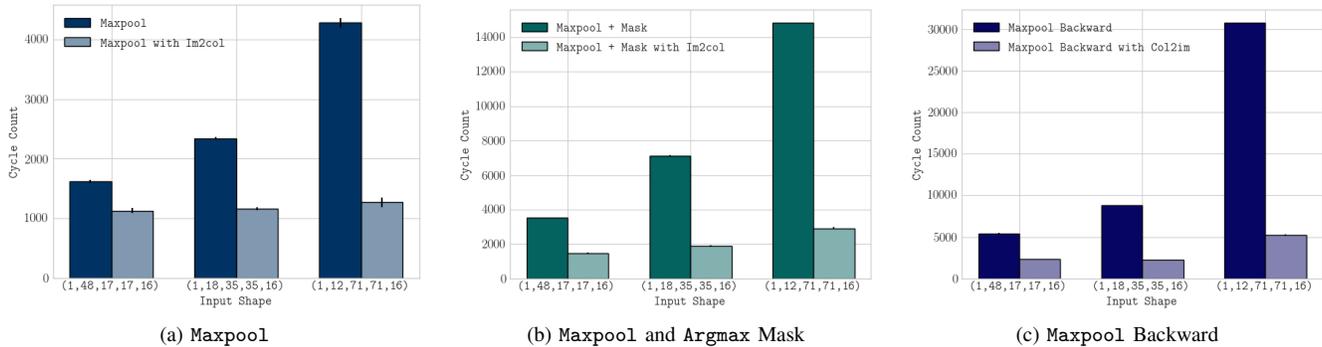


Fig. 7. Comparison of Maxpool implementations with and without Im2Col and Col2Im instructions. The graphs show the cycle count in the Ascend 910 chip by the size of the input. The input sizes are from InceptionV3. All tests use a kernel size of (3,3) and a stride of (2,2) with no padding.

TABLE I
MAXPOOL INPUT SIZES IN CNNs

CNN	Input 1	Input 2	Input 3	Input 4
InceptionV3	147,147,64	71,71,192	35,35,288	17,17,768
Xception	147,147,128	74,74,256	37,37,728	19,19,1024
Resnet50	112,112,64	-	-	-
VGG16	224,224,64	112,112,128	56,56,256	28,28,512

578 TVM, they are declared and manually added to the code as
 579 custom intrinsics through TVM’s `decl_tensor_intrin`
 580 function. These intrinsics act in TVM’s DSL as an inline
 581 assembly section in a C source. Instead of implementing a
 582 single instruction call, the custom intrinsics were defined to
 583 issue instructions multiple times. By also using the repetition
 584 parameters, they can operate on a full tile of the input.

585 A. InceptionV3 Comparison

586 Table I shows multiple CNNs and the input sizes of four
 587 of their Maxpool layers. The inputs are shown in the HWC
 588 layout and they were gathered on the Keras framework [30].
 589 All configurations use a kernel size of (3,3) and a stride
 590 of (2,2), except for VGG16 [31], which has a kernel size
 591 and stride of (2,2). To test the implementations of Maxpool,
 592 three configurations were selected from InceptionV3 [32]
 593 (highlighted in bold). No padding is used in them, however,
 594 it is also possible to add padding during the Im2Col load, as
 595 the other CNNs would require. Given AKG’s current limited
 596 support for the Im2Col and Col2Im, these configurations were
 597 chosen to display the effects of different input sizes while
 598 using the most common parameters of kernel and stride.

599 The graphs in Figure 7 show the cycle count of the selected
 600 Maxpool configurations in the NC_1HWC_0 layout. Figure 7a
 601 shows both Maxpool forward implementations. The step of
 602 saving the Argmax mask is added in Figure 7b. This step adds
 603 to the computation, as shown by the different ranges in the
 604 graphs. For the evaluation in Figure 7b, AKG’s polyhedral
 605 framework schedules the computations, as it can better handle
 606 computations with multiple outputs of different shapes. Lastly,
 607 Maxpool backward is evaluated in Figure 7c. In the largest

input, the accelerated implementations achieve speedups of
 3.2x, 5x, and 5.8x on the graphs in Figure 7, respectively. The
 best improvement is on Maxpool backward. Its large speedup
 is expected, given the scattered access pattern of its merge step
 and how Col2Im can be used without any extra computations.

B. Stride Tests

608 This experiment investigates further different Maxpool for-
 609 ward implementations, and their interaction with the stride
 610 parameter, as shown in Figure 8. The stride size changes
 611 the amount of duplicated elements in Im2col. The kernel
 612 size was set at a constant size of (3,3). Given this kernel
 613 size, there is no duplication of data for the (3,3) stride,
 614 and the maximum duplication occurs for the (1,1) stride.
 615 In this experiment, Maxpool and Maxpool with Im2col are
 616 the same implementations shown in Figure 7a. The input’s
 617 height and width increase in steps of two until the tiling
 618 threshold is reached, where this threshold is the maximum size
 619 before tiling is required. Bigger sizes would need individual
 620 tiling parameters and would trigger parallelization between AI
 621 Cores, which is out of the scope of this experiment. Moreover,
 622 dimensions N and C_1 are set to 1 so that only one AI Core
 623 is utilized.

624 In the Maxpool with expansion implementation, reg-
 625 ular vector instructions — instead of Im2Col instructions
 626 — transform the input to the Im2Col output shape. This
 627 transformation happens when the input is already in the
 628 Unified buffer, before computing Maxpool. Maxpool with
 629 Im2col and Maxpool with expansion achieve superior per-
 630 formance in Figures 8b and 8c. These graphs confirm that
 631 the Im2col memory layout allows more efficient usage of
 632 the Vector Unit, producing speedups that compensate for the
 633 overhead of transforming the data. Maxpool with Im2col
 634 has the best performance in comparison to Maxpool with
 635 expansion due to Im2col occurring while the data is loaded
 636 into the Unified buffer, rather than in a separate step.

637 Figure 8a shows different results for a stride of (1,1). With
 638 this parameter, elements in consecutive patches of the original
 639 input appear consecutively in memory. This allows the `vmax`
 640 instruction to improve its use of the Vector Unit, combining the
 641 mask register set with all 128 elements and its repeat parameter
 642

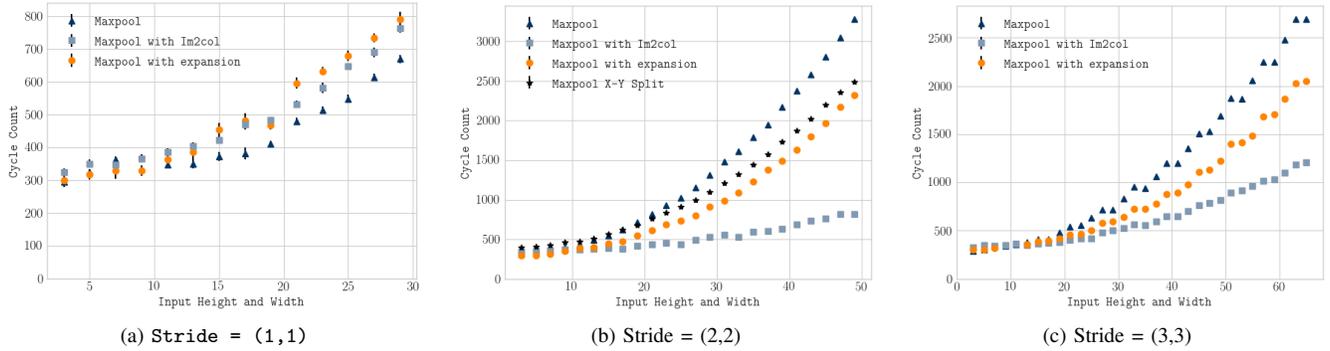


Fig. 8. Comparison of different Maxpool implementations. The graphs show the cycle count in the Ascend 910 board and the height and width of the input. In all tests, the N and C_1 sizes are 1, kernel size is (3,3), with no padding. The x-axis goes up to the tiling threshold. An additional implementation of the X-Y split is shown for the stride of (2,2).

648 to compute the max between the (O_w, C_0) dimensions of the
 649 input and the initialized output. By also having no overhead to
 650 transform the data, and no data duplication, the direct Maxpool
 651 implementation is the fastest in this case.

652 Pooling can also be implemented with an X-Y split by first
 653 calculating the reduction function on the width and then on the
 654 height of each patch. As a result, the first reduction is reused
 655 while computing the second. Lai *et al.* [7] use the X-Y split as
 656 a performant alternative to direct pooling. In their work, the
 657 (undesirable) intermediate results are avoided by computing
 658 the result in-place. In TVM, all computations generate a
 659 new tensor, and thus the in-place approach is not possible.
 660 However, this experiment increases input sizes only until the
 661 tiling threshold is reached, therefore avoiding extra tiling steps
 662 needed because of the increase in memory use. Figure 8b
 663 shows the performance of a TVM version of the X-Y split
 664 (with intermediate results) compared to the other Maxpool
 665 implementations. Even though the X-Y split has a lower
 666 computational cost, it underperforms other implementations
 667 that use Im2Col because it does not overcome the scattered
 668 memory problems of pooling.

669 VII. RELATED WORK

670 Convolutional layers have been the focus of extensive liter-
 671 ature in optimizing CNN layers because they are responsible
 672 for most of the computation time of CNNs. Other layers such
 673 as pooling receive less attention, but when left unoptimized,
 674 they can be obstacles that lead to slowdowns in CNNs [1].

675 **FPGA implementations for CNNs.** In their implementa-
 676 tion of CNN layers for OpenCL-based FPGA accelerators,
 677 Suda *et al.* propose to unroll pooling at the hardware level
 678 so that multiple outputs are computed in a single cycle [5].
 679 However, their optimizer chooses an unrolling factor of 1 for
 680 the CNNs evaluated, which is equal to no unrolling. Given
 681 an (FPGA, CNN) pair, Sharma *et al.* automatically synthe-
 682 size a CNN accelerator where the computation of pooling
 683 modules overlaps with convolution modules. This overlap
 684 is used to hide latency and to take advantage of the fact
 685 that a pooling layer usually follows convolutional layers [6].
 686 Sharma *et al.* do not consider the backward operators used

687 in training. In contrast to these pooling-specific hardware
 688 solutions, Im2col/Col2Im based pooling in DaVinci leverages
 689 a general-purpose vector unit and the Im2Col and Col2Im
 690 instructions, which are primarily designed for convolution. The
 691 improvements to the pooling layer afforded by Im2col/Col2im
 692 could be combined with fusion in DaVinci, but this is not yet
 693 supported.

694 **Kernel acceleration for CNNs.** LightNet is a Matlab-based
 695 framework for Deep Learning [33]. Its Maxpool implementa-
 696 tion uses Im2col to transform pooled regions into vectors to
 697 benefit from vector instructions. Their proposition is similar to
 698 the Im2col based forward pooling, however, no performance
 699 results are presented to justify their implementation. CMSIS-
 700 NN is a collection of efficient neural network layers targeting
 701 IoT edge devices that uses X-Y splitting for pooling [7].
 702 However, the results in Figure 8b show that the X-Y split is
 703 not the best alternative for DaVinci. Additionally, CMSIS-NN
 704 does not consider backward operators because its target edge
 705 devices only perform inference. The Im2col/Col2im based
 706 pooling accelerates both inference and training devices, as
 707 DaVinci edge chips also feature Im2Col instructions.

708 Li *et al.* use two optimizations for pooling [34]. First, the
 709 use of the $CHWN$ layout instead of $NCHW$ to prevent un-
 710 coalesced strided memory accesses caused by HW as the
 711 innermost dimensions. Second, the reduction of the off-chip
 712 memory requests by tuning the number of outputs calculated
 713 by each thread during pooling. The memory layout used in
 714 DaVinci (NC_1HWC_0) is a variant of the $NCHW$ layout.
 715 However, the Im2col-based pooling transforms this layout
 716 into $NC_1K_hK_wO_hO_wC_0$, where the accesses can also be
 717 performed consecutively in memory, thus resulting in the
 718 performance speedups shown in Section VI. The outer loops
 719 are automatically parallelized in DaVinci among the available
 720 AI Cores, where each core calculates a share of the output.

721 Suita *et al.* focus on fusing convolution with pooling in
 722 GPUs [8]. They only consider Avgpool because it can be
 723 mapped to convolution where the kernel's weights are equal
 724 to $1/(K_h * K_w)$, and then further fused with its preceding
 725 convolution. As a result, the Im2col transformation can also

726 be used to implement the fused convolution-pooling. However,
 727 CNNs tend to use Maxpool, which cannot be fused in the same
 728 way.

729 VIII. CONCLUSION

730 This work presents DaVinci's Im2Col and Col2Im instruc-
 731 tions. It is shown that they can be used to implement not
 732 only convolution, targeting the Cube Unit, but also pooling,
 733 targeting memory layout improvements for the Vector Unit.
 734 Accelerated Im2col/Col2im based implementations are de-
 735 scribed for the forward and backward operators of Maxpool
 736 and Avgpool. An experimental evaluation was run on the
 737 Ascend 910 chip with the parameters and three input sizes
 738 used in InceptionV3. The results show speedups of up to 5.8x
 739 for the accelerated Maxpool implementations, compared to
 740 baselines that do not use the Im2Col and Col2Im instructions.
 741 Although the stride parameter can impact the Im2col and
 742 Col2im operations drastically, the proposed acceleration ap-
 743 proach achieved improved performance for all but (1, 1) stride.
 744 The Im2col/Col2im based pooling also proves superior to other
 745 strategies of optimization, such as the X-Y split. Further work
 746 could evaluate the proposed approach in other architectures,
 747 and also consider the fusion techniques described by Suita *et*
 748 *al.* [8] to execute Avgpool together with convolution as matrix
 749 multiplication in the Cube Unit.

750 REFERENCES

- 751 [1] D. Li, X. Wang, and D. Kong, "DeepRebirth: accelerating deep neural
 752 network execution on mobile devices," in *AAAI Conference on Artificial*
 753 *Intelligence*, New Orleans, LA, USA, 2018, pp. 2322–2330.
- 754 [2] H. Liao, J. Tu, J. Xia, and X. Zhou, "DaVinci: A scalable architecture for
 755 neural network computing," in *Hot Chips Symposium (HCS)*, Cupertino,
 756 CA, USA, 2019, pp. 1–44.
- 757 [3] A. Khaled, A. F. Atiya, and A. H. Abdel-Gawad, "Applying fast matrix
 758 multiplication to neural networks," in *Symposium on Applied Computing*
 759 *(SAC)*, Brno, Czech Republic, 2020, p. 1034–1037.
- 760 [4] Y. Guan *et al.*, "FP-DNN: an automated framework for mapping deep
 761 neural networks onto FPGAs with RTL-HLS hybrid templates," in *Field-*
 762 *Programmable Custom Computing Machines (FCCM)*, Napa, CA, USA,
 763 2017, pp. 152–159.
- 764 [5] N. Suda *et al.*, "Throughput-optimized OpenCL-based FPGA accelerator
 765 for large-scale convolutional neural networks," in *Field-Programmable*
 766 *Gate Arrays (FPGA)*, Monterey, California, USA, 2016, p. 16–25.
- 767 [6] H. Sharma *et al.*, "From high-level deep neural models to FPGAs," in
 768 *Microarchitecture (MICRO)*, Taipei, Taiwan, 2016, pp. 1–12.
- 769 [7] L. Lai, N. Suda, and V. Chandra, "CMSIS-NN: efficient neural network
 770 kernels for Arm Cortex-M CPUs," *CoRR*, vol. abs/1801.06601, 2018.
- 771 [8] S. Suita *et al.*, "Efficient convolution pooling on the gpu," *Journal of*
 772 *Parallel and Distributed Computing*, vol. 138, pp. 222 – 229, 2020.
- 773 [9] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning
 774 applied to document recognition," *Proceedings of the IEEE*, vol. 86,
 775 no. 11, pp. 2278–2324, 1998.
- 776 [10] A. Vasudevan, A. Anderson, and D. Gregg, "Parallel multi channel
 777 convolution using general matrix multiplication," in *Application-specific*
 778 *Systems, Architectures and Processors (ASAP)*, Seattle, WA, USA, 2017,
 779 pp. 19–24.
- 780 [11] K. Chellapilla, S. Puri, and P. Simard, "High performance convolutional
 781 neural networks for document processing," in *Tenth International Work-*
 782 *shop on Frontiers in Handwriting Recognition*, La Baule (France), 2006.
- [12] Z. Xianyi, W. Qian, and Z. Yunquan, "Model-driven level 3 BLAS
 performance optimization on loongson 3A processor," in *International*
Conference on Parallel and Distributed Systems (ICPADS), Singapore,
 2012, p. 684–691.
- [13] R. C. Whaley and J. J. Dongarra, "Automatically tuned linear algebra
 software," in *International Conference for High Performance Comput-*
ing, Networking, Storage and Analysis (SC), San Jose, CA, USA, 1998,
 p. 1–27.
- [14] G. Guennebaud, B. Jacob *et al.* (2010) Eigen v3. [Online]. Available:
<http://eigen.tuxfamily.org>
- [15] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, *Learning Repre-*
sentations by Back-Propagating Errors. Cambridge, MA, USA: MIT
 Press, 1988, p. 696–699.
- [16] D. A. Forsyth and J. Ponce, *Computer Vision: A Modern Approach*.
 Prentice Hall Professional Technical Reference, 2002.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image
 recognition," in *Computer Vision and Pattern Recognition (CVPR)*, Las
 Vegas, NV, USA, 2016, pp. 770–778.
- [18] C. Szegedy *et al.*, "Going deeper with convolutions," in *Computer Vision*
and Pattern Recognition (CVPR), Boston, MA, USA, 2015, pp. 1–9.
- [19] F. Chollet, "Xception: Deep learning with depthwise separable convolu-
 tions," in *Computer Vision and Pattern Recognition (CVPR)*, Honolulu,
 HI, USA, 2017, pp. 1800–1807.
- [20] F. Chollet, *Deep Learning with Python*, 1st ed. Manning Publications
 Co., 2017.
- [21] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind,
 "Automatic differentiation in machine learning: A survey," *J. Mach.*
Learn. Res., vol. 18, no. 1, p. 5595–5637, 2017.
- [22] D. Scherer, A. Müller, and S. Behnke, "Evaluation of pooling opera-
 tions in convolutional architectures for object recognition," in *Internat-*
ional Conference on Artificial Neural Networks (ICANN), Thessaloniki,
 Greece, 2010, p. 92–101.
- [23] H. T. Kung and C. E. Leiserson, "Systolic arrays for (VLSI)." Carnegie-
 Mellon Univ Pittsburgh Pa Dept Of Computer Science, Tech. Rep., 1978.
- [24] N. P. Jouppi *et al.*, "In-datacenter performance analysis of a tensor
 processing unit," in *International Symposium on Computer Architecture*
(ISCA), Toronto, ON, Canada, 2017, p. 1–12.
- [25] B. Jacob, S. Ng, and D. Wang, *Memory Systems: Cache, DRAM, Disk*.
 San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2007.
- [26] T. Chen *et al.*, "TVM: An automated end-to-end optimizing compiler
 for deep learning," in *Operating Systems Design and Implementation*
(OSDI), Carlsbad, CA, USA, 2018, p. 579–594.
- [27] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Ama-
 rasinghe, "Halide: A language and compiler for optimizing parallelism,
 locality, and recomputation in image processing pipelines," in *Program-*
ming Language Design and Implementation (PLDI). New York, NY,
 USA: Association for Computing Machinery, 2013, p. 519–530.
- [28] M. E. Wolf and M. S. Lam, "A loop transformation theory and an
 algorithm to maximize parallelism," *IEEE Trans. Parallel Distrib. Syst.*,
 vol. 2, no. 4, p. 452–471, 1991.
- [29] R. T. Mullaipudi, A. Adams, D. Sharlet, J. Ragan-Kelley, and K. Fa-
 tahalian, "Automatically scheduling halide image processing pipelines,"
ACM Trans. Graph., vol. 35, no. 4, 2016.
- [30] F. Chollet *et al.*, "Keras," <https://keras.io>, 2015.
- [31] K. Simonyan and A. Zisserman, "Very deep convolutional networks for
 large-scale image recognition," in *International Conference on Learning*
Representation (ICLR), San Diego, CA, USA, 2015.
- [32] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking
 the inception architecture for computer vision," in *Computer Vision and*
Pattern Recognition (CVPR), Las Vegas, NV, USA, 2016, pp. 2818–
 2826.
- [33] C. Ye, C. Zhao, Y. Yang, C. Fermüller, and Y. Aloimonos, "LightNet:
 A versatile, standalone matlab-based environment for deep learning," in
Multimedia (MM), Amsterdam, The Netherlands, 2016, p. 1156–1159.
- [34] C. Li, Y. Yang, M. Feng, S. Chakradhar, and H. Zhou, "Optimizing
 memory efficiency for deep convolutional neural networks on GPUs," in
International Conference for High Performance Computing, Networking,
Storage and Analysis (SC), Salt Lake City, Utah, 2016.