

An Efficient Method for Generic DSP Implementation of Dilated Convolution

Harinarayanan.E.V and Sachin Ghanekar
Cadence Design Systems, Pune, India

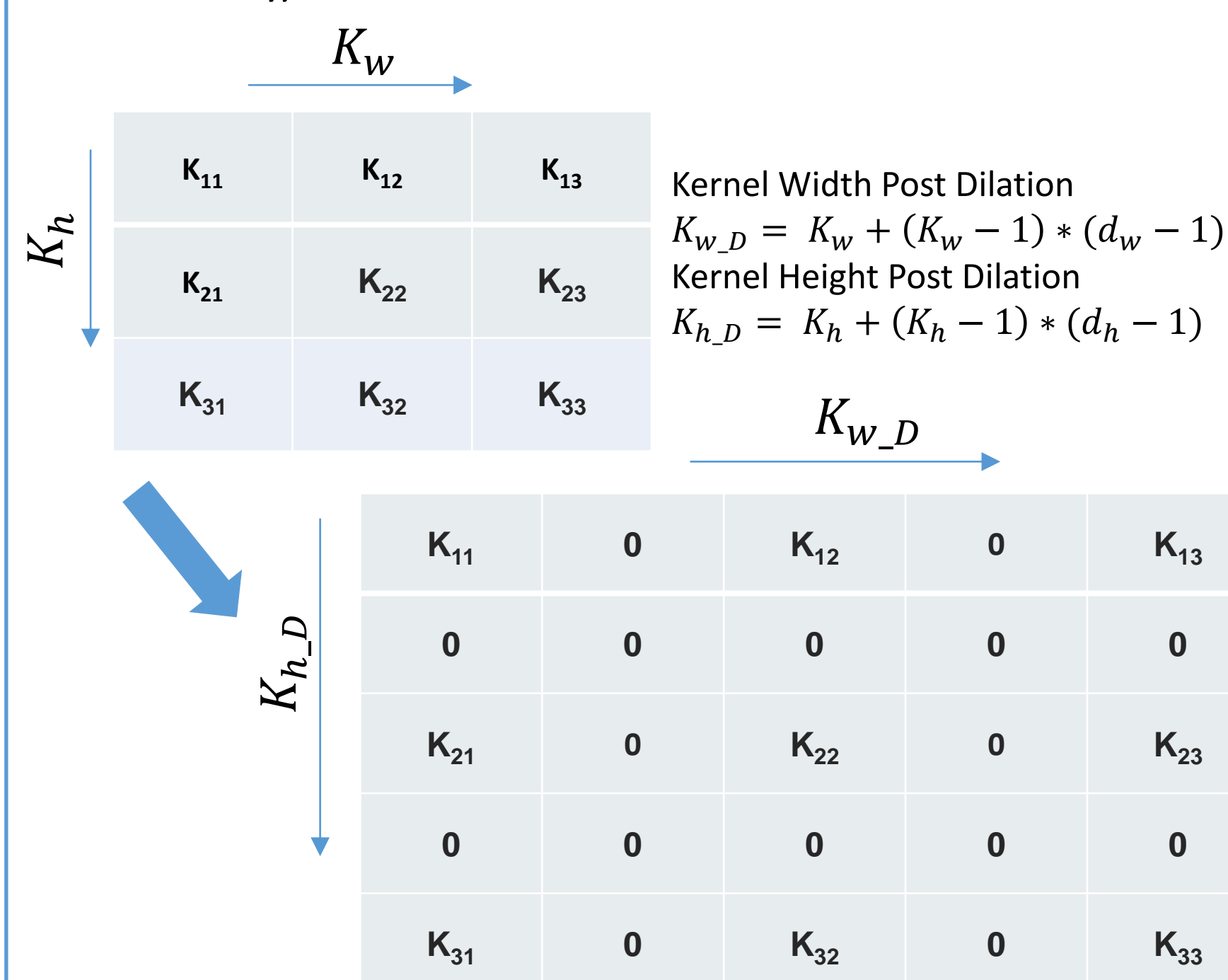
Introduction

Dilated convolution has an inherent property of capturing wider context in an image and long-term temporal characteristics in an audio signal. We propose a new scheme that allows efficient/generic implementation of 2D Dilated convolution and stride on typical DSPs where the instruction sets are well tuned for standard 1D and 2D filtering and convolution operations.

In this proposal an existing flexible and efficient standard 2D convolution implementation with stride support forms the basic building block to implement Dilated convolution with stride. Thereby, 2D-dilated convolution is equivalently represented as several smaller convolutions with appropriate matrix slicing and re-ordering.

Dilation & Stride

2D Dilation in literature referred to as *atrous* convolution or convolution with holes. Dilation introduces 'zeros' in the kernel matrix of the convolution. An example of pre/post kernel dilation with a factor 2 in height d_h and width d_w shown below,



Matrices post dilated convolution are generally strided to reduce the dimensions of the output. Stride simply-put is skipping the matrix values by a factor. Stride and Dilation in mathematical form for 2D matrix is as below:

$$Y(x_3, y_3) = \sum_{0 \leq x_2 < K_h} \sum_{0 \leq y_2 < K_w} I \left(\begin{matrix} x_2 * d_h + x_3 * s_h \\ y_2 * d_w + y_3 * s_w \end{matrix} \right) * F(x_2, y_2)$$

Where, I the input matrix, F kernel matrix and s_w/s_h are the stride factors in width and height dimension, $0 \leq y_3 < [(I_w - K_{w,D})/s_w] + 1$ and $0 \leq x_3 < [(I_h - K_{h,D})/s_h] + 1$

The input matrices are assumed to be appropriately zero padded, if needed, such that input width and input height are always gr. than equal to kernel dimension i.e., $I_w \geq K_{w,D}$ and $I_h \geq K_{h,D}$.

Re-ordering stride output

$$Y \left(\frac{h_{offset}(n_1) + i_1 * R_h}{s_h}, \frac{w_{offset}(n_2) + i_2 * R_w}{s_w} \right) = Y_{n_1 n_2 d_h d_w s_h s_w}(i_1, i_2)$$

$$n_1 = 0, 1, 2, \dots, d_h - 1$$

$$n_2 = 0, 1, 2, \dots, d_w - 1$$

$$h_{offset}(n_1) = n_1 + x_{n_{min}} * d_h$$

$$w_{offset}(n_2) = n_2 + y_{n_{min}} * d_w$$

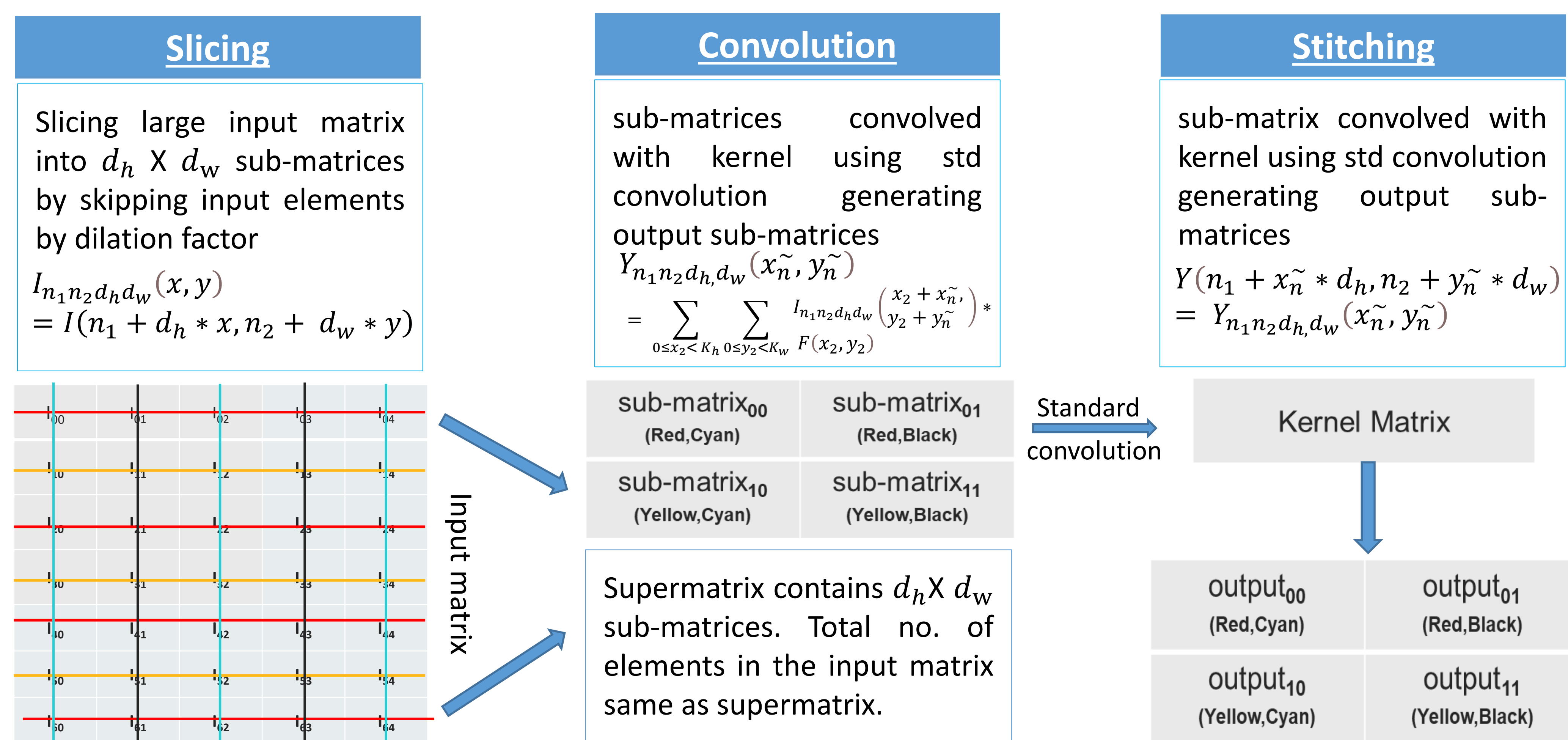
$$R_h = \frac{s_h}{GCD(d_h, s_h)} * d_h; R_w = \frac{s_w}{GCD(d_w, s_w)} * d_w$$

A Novel scheme for Dilation and Stride

Dilated convolution in comparison to non-dilation increases the computational complexity by the order of the product of dilation factors in height and width dimension using *atrous* method. We propose a joint solution for dilation and stride. As a special case, the complexity is immune to the dilation factor when stride is unity.

Dilation Scheme

Implementation of Dilated convolution split into three steps: a) Input Slicing b) Standard convolution c) Output stitching. An existing standard convolution forms the central computation block and input-slicing/output-stitching are the memory alignment processes. Below an illustration of a 7X5 matrix for dilation factor of 2 in height and width.



Stride Scheme

The skipped values of convolved output by a pre-defined factor is given by stride. The scheme for stride is proposed in the dilation framework as explained above i.e., slicing/convolution/stitching

Problem Statement

From the 'stitching' step the output values of interest after applying stride can be written as,

$$Y(s_h * x, s_w * y) = Y(n_1 + x_{n_1} * d_h, n_2 + y_{n_1} * d_w)$$

Where, s_h and s_w are stride in height and width dimension. In other words, matrix values with integer coordinates of $\langle x, y \rangle$ are values of interest as shown below.

$$\langle x, y \rangle = \left\langle \frac{n_1 + x_{n_1} * d_h}{s_h}, \frac{n_2 + y_{n_1} * d_w}{s_w} \right\rangle$$

Alternate Statement

S1: For a given offset pair $\langle n_1, n_2 \rangle$ find the minimum value of coordinates $\langle x_{n_{min}}, y_{n_{min}} \rangle$ say, $\langle x_{n_{min}}, y_{n_{min}} \rangle$

S2: From the initial values $\langle x_{n_{min}}, y_{n_{min}} \rangle$ for a given $\langle n_1, n_2 \rangle$ find successive values of $\langle x_{n_1}, y_{n_1} \rangle$

Solution

Postulate: Assume S1 is true i.e., for a given $\langle n_1, n_2 \rangle$ there exist a $\langle x_{n_{min}}, y_{n_{min}} \rangle$. Finding solution for S2 i.e., to find the successive values of $\langle x_{n_1}, y_{n_1} \rangle$ after the initial value $\langle x_{n_{min}}, y_{n_{min}} \rangle$

Inference: Let the next successive value of $\langle x_{n_1}, y_{n_1} \rangle$ after $\langle x_{n_{min}}, y_{n_{min}} \rangle$ be $\langle x_{n_{min}} + \Delta_x, y_{n_{min}} + \Delta_y \rangle$. Inserting these values.

$$\langle x, y \rangle = \left\langle \frac{n_1 + (x_{n_{min}} * d_h)}{s_h} + \frac{\Delta_x * d_h}{s_h}, \frac{n_2 + (y_{n_{min}} * d_w)}{s_w} + \frac{\Delta_y * d_w}{s_w} \right\rangle$$

The minimum value of Δ_x & Δ_y to contribute an integer value $\langle x, y \rangle$ is

$$\Delta_x, \Delta_y = \frac{s_h}{GCD(s_h, d_h)}, \frac{s_w}{GCD(s_w, d_w)}$$

Therefore $\langle \Delta_x, \Delta_y \rangle$ is the periodic pattern for chosen value of dilation and stride. Extending this periodic property, it is sufficient to check the first $\langle \Delta_x, \Delta_y \rangle$ elements to test the validity of statement S1.

Highlights:

- A sub-matrix $\langle n_1, n_2 \rangle$ not satisfying S1 => no participation in convolution
- Δ_x & Δ_y are the modified stride values for sub-matrix convolution
- Upon re-ordering $\frac{d_h}{GCD(s_h, d_h)}, \frac{d_w}{GCD(s_w, d_w)}$ is the output stride

Result

The proposed method of decomposition is compared against *atrous* method. The method implemented on Cadence's Tensilica HiFi5 processor with NN extension simulator assuming zero memory wait states. Input/kernel data format (N,H,W,C). An improvement of 30X (cycles) can be observed for dilation factor of 16. Scratch memory reduction for the proposed method observed. The implementation is available on Cadence's NN HiFi5 [GitHub link](#).

| Dilation Factor | Decomposition Method | | Zero Insertion (ZI) Method | |
|-----------------|----------------------------|---------------------|----------------------------|---------------------|
| | Cycles (X10 ⁶) | Scratch Memory (KB) | Cycles (X10 ⁶) | Scratch Memory (KB) |
| (Stride=2) | | | | |
| 2 | 1.94 | 6.14 | 4.89 | 20.20 |
| 4 | 2.01 | 3.14 | 10.01 | 36.33 |
| 8 | 2.14 | 1.64 | 26.54 | 68.58 |
| 16 | 2.32 | 0.89 | 69.70 | 133.08 |

Computational gain for different Dilation and Stride values published in paper