

Future Design Systems	FDS-TD-2020-09-003

# High-Level Synthesizable Deep Learning Routines with Python/PyTorch Interfaces

Version 0 Revision 6

Dec. 17, 2020 (July 23, 2020)

Future Design Systems, Inc.

[www.future-ds.com](http://www.future-ds.com) / [contact@future-ds.com](mailto:contact@future-ds.com)

**Copyright © 2020 Future Design Systems, Inc.**

## Abstract

This document addresses Deep Learning Routines that are High-Level synthesizable. In addition to the C routines, Python Interfaces for PyTorch are included in order to verify functionality of the routines along with PyTorch deep learning framework.

## Table of Contents

Copyright © 2020 Future Design Systems, Inc. ....	1
Abstract .....	1
Table of Contents .....	1
1 Overview .....	4
1.1 Verification flow .....	4
2 License .....	5
3 Convention and background.....	5
3.1 Naming convention.....	5
3.2 Multi-dimensional array .....	6
3.3 Python wrapper .....	6
3.4 PyTorch wrapper .....	6
3.5 Hardware interface .....	6
3.6 Software interface .....	7
4 Activation .....	8
4.1 ReLU .....	8
4.1.1 Python wrapper .....	9
4.1.2 PyTorch wrapper .....	9
4.2 Leaky ReLU .....	9
4.2.1 Python wrapper .....	10
4.2.2 PyTorch wrapper .....	10

4.3 Hyperbolic tangent .....	10
4.3.1 Python wrapper .....	11
4.3.2 PyTorch wrapper .....	11
4.4 Sigmoid .....	11
4.4.1 Python wrapper .....	12
4.4.2 PyTorch wrapper .....	12
5 Concatenation .....	13
5.1 Concat 2D .....	13
5.1.1 Python wrapper .....	14
5.1.2 PyTorch wrapper .....	14
6 Convolution .....	14
6.1 Convolution2d .....	14
6.1.1 Python wrapper .....	16
6.1.2 PyTorch wrapper .....	17
7 Deconvolution (Transposed convolution) .....	17
7.1 deconvolution2d .....	17
7.1.1 Python wrapper .....	18
7.1.2 PyTorch wrapper .....	19
8 Linear (Fully connected) .....	19
8.1 Linear 1D.....	19
8.1.1 Python wrapper .....	20
8.1.2 PyTorch wrapper .....	20
8.2 Linear ND .....	21
8.2.1 Python wrapper .....	22
8.2.2 PyTorch wrapper .....	22
9 Normalization .....	22
9.1 Batch normalization 2D .....	22
9.1.1 Python wrapper .....	24
9.1.2 PyTorch wrapper .....	24
10 Pooling .....	24
10.1 Pooling2dMax .....	24
10.1.1 Python wrapper .....	25
10.1.2 PyTorch wrapper .....	26
10.2 Pooling2dAvg .....	26
10.2.1 Python wrapper .....	27
10.2.2 PyTorch wrapper .....	27
11 Project: LeNet-5.....	28
11.1 LeNet-5.....	28
11.1.1 LeNet-5 C++ version .....	28
11.1.2 LeNet-5 C PyTorch version.....	32
12 Troubleshooting.....	32
13 Acknowledgement .....	33
14 References .....	33
Wish list.....	33
Index .....	33
Revision history .....	34

Future Design Systems	FDS-TD-2020-09-003

Future Design Systems	FDS-TD-2020-09-003

## 1 Overview

**DLR** (*Deep Learning Routines*) as a part of **DPU** (*Deep Learning Processing Unit*) is a collection of high-level synthesizable C routines for deep learning inference network. **DLIP** (*Deep Learning Intellectual Property*) is HDL (Hardware Description Language) version of DLR.

DLR contain followings<sup>1</sup>:

- ✧ Convolution layer (strict convolution, matrix-multiplication-based convolution)
- ✧ Activation layer (ReLU, Leaky ReLU, hyperbolic tangent, sigmoid)
- ✧ Pooling layer (max, average)
- ✧ Fully connected layer (linear layer)
- ✧ De-convolution layer (transposed convolution layer)
- ✧ Batch normalization layer
- ✧ Concatenation layer

The routines have following highlights:

- ✧ Fully synthesizable C code
- ✧ Highly parameterized to be adopted wide range of usages
- ✧ C, Python, PyTorch, Verilog test-benches
- ✧ FPGA verified using Future Design Systems' CON-FMC

It should be noted that the routines are not optimized to get a higher performance since the routines are for hardware implementation not for computation. In addition to this the routines are only for inference not for training.

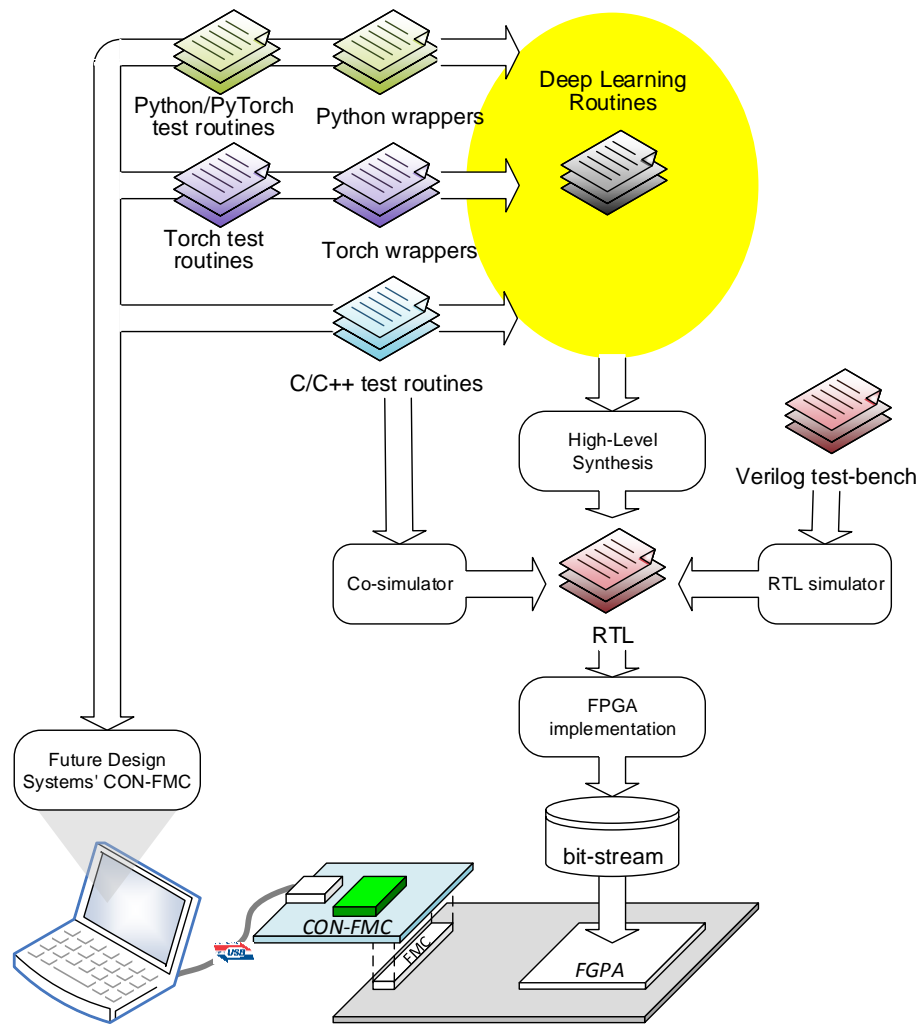
### 1.1 Verification flow

Figure 1 shows an overall verification flow, in which following interfaces are supported.

- Standard C/C++
- PyTorch C++ Front-End
- PyTorch

---

<sup>1</sup> Some are under development and more will be added.



**Figure 1: Verification flow**

## 2 License

DLR (Deep Learning Routines) and its associated materials are licensed with the 2-clause BSD license to make the program and library useful in open and closed source products independent of their licensing scheme.

Copyright © 2020 by Future Design Systems, Inc.

<http://www.future-ds.com>

[contact@future-ds.com](mailto:contact@future-ds.com)

## 3 Convention and background

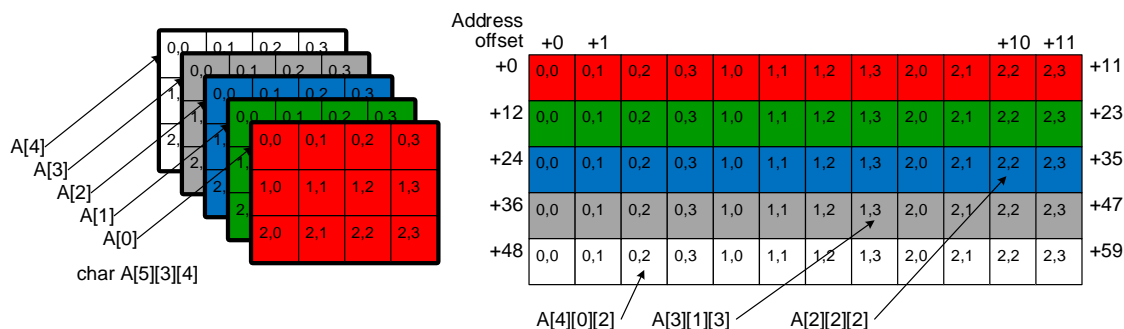
### 3.1 Naming convention

This routine uses following coding conventions.

- ✧ 4 spaces for indentation, not tab
- ✧ Variables uses underscores separating words.
  - char item\_delimiter;
- ✧ Functions starts with upper case character and follows camel-case.
  - void Convolution2dWrapper(int stride);
- ✧ Macros uses all capital letter with underscores separating words.
  - #define MAX\_NUM 128
- ✧ Temporary variable starts underscore or 't\_' prefix.
  - int t\_tmp;
- ✧ Output comes first in the function argument list.
- ✧ Python wrapper functions follow format of DRL of C.
- ✧ PyTorch wrapper functions follow format of PyTorch nn.functional.

### 3.2 Multi-dimensional array

This routine uses row-major contiguous array layout, which is the same as that of the standard C language.



**Figure 2: Multi-dimensional array (one-byte data case)**

### 3.3 Python wrapper

Python wrapper uses Ctypes (a foreign function library for Python) to call C routine.

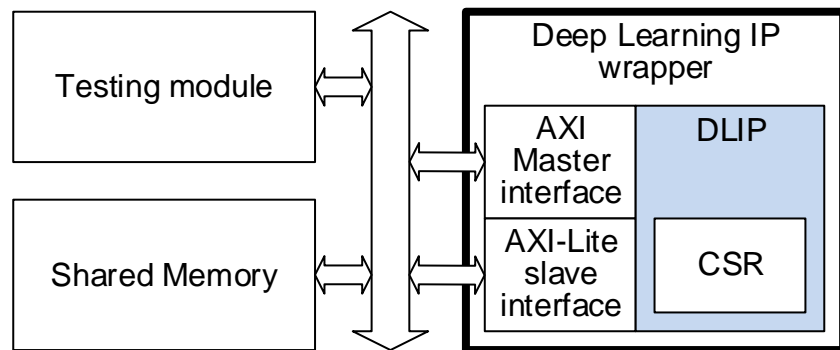
DLIP (Deep Learning IP) Python wrapper uses Python NumPy array to carry array data. As a result, PyTorch tensor should be converted to NumPy array in order to be called from PyTorch.

### 3.4 PyTorch wrapper

PyTorch wrapper supports torch.nn.functional that is stateless functional version.

### 3.5 Hardware interface

DLIP (Deep Learning IP) has two AMB AXI bus interfaces; one is for accessing internal registers and the other is for access shared memory.



**Figure 3: Deep Learning IP HW structure including test bench**

### 3.6 Software interface

Following explains the structure of CSR (Control and Status Register) in the DLIP.

```
// ctl
// 0x00 : Control signals
//   bit 0 - ap_start (Read/Write/COH)
//   bit 1 - ap_done (Read/COR)
//   bit 2 - ap_idle (Read)
//   bit 3 - ap_ready (Read)
//   bit 7 - auto_restart (Read/Write)
//   others - reserved
// 0x04 : Global Interrupt Enable Register
//   bit 0 - Global Interrupt Enable (Read/Write)
//   others - reserved
// 0x08 : IP Interrupt Enable Register (Read/Write)
//   bit 0 - Channel 0 (ap_done)
//   bit 1 - Channel 1 (ap_ready)
//   others - reserved
// 0x0c : IP Interrupt Status Register (Read/TOW)
//   bit 0 - Channel 0 (ap_done)
//   bit 1 - Channel 1 (ap_ready)
//   others - reserved
// 0x10 : Data signal of shared_mem
//   bit 31~0 - shared_mem[31:0] (Read/Write)
// 0x14 : reserved
// (SC = Self Clear, COR = Clear on Read, TOW = Toggle on Write, COH = Clear on Handshake)
```

Following shows an example to control DLIP through CSR.

Wait for ready before start	Let DLIP runs and wait for completion
#define WAIT_FOR_READY {\n  int ap_idle, ap_idle_r;\n  unsigned int ap_addr;\n}	#define GO_AND_WAIT_COMPLETE {\n  int ap_done, ap_done_r;\n  int ap_start, ap_data;\n}

Future Design Systems	FDS-TD-2020-09-003

<pre> ap_addr = ADDR_CSR;\ while (1) {\     MEM_READ(ap_addr, ap_idle_r);\     ap_idle = (ap_idle_r &gt;&gt; 2) &amp;&amp; 0x1;\     if (ap_idle) break;\ } </pre>	<pre> unsigned int ap_addr;\ ap_addr = ADDR_CSR;\ ap_data = 0x1;\ MEM_WRITE(ap_addr, ap_data);\ while (1) {\     MEM_READ(ap_addr, ap_done_r);\     ap_done = (ap_done_r &gt;&gt; 1) &amp;&amp; 0x1;\     if (ap_done) break;\ } </pre>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

## 4 Activation

### 4.1 ReLU

‘ActivationReLU()’ applies Rectified Linier Unit over an input data.

<pre> #include "activation_relu.hpp"  template&lt;class TYPE=float&gt; void ActivationReLu (     TYPE      *out_data     , const TYPE *in_data     , const uint32_t size     , const uint16_t channel     #if !defined(__SYNTHESIS__)     , const int    rigor=0     , const int    verbose=0     #endif ); </pre>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Templates:

- ✧ TYPE: data type of out\_data, in\_data
  - int, float, double, and so on<sup>2</sup>.

Arguments:

- ✧ out\_data: resultant data in contiguous memory and can be any dimension
  - channel x size x size
- ✧ in\_data: input array in contiguous memory and can be any dimension
  - channel x size x size
- ✧ size: the number of elements per channel
- ✧ channel: the number of channels
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- ‘out\_data’ and ‘in\_array’ must have ‘channel x size’ elements.

<sup>2</sup> Fixed-point data type will be supported.



Future Design Systems	FDS-TD-2020-09-003

#### 4.1.1 Python wrapper

‘ActivationReLU()’ wrapper gets NumPy arguments for array and ‘out\_data’ carries calculated result. It returns ‘True’ on success or ‘False’ on failure.

```
ActivationReLU( out_data # any dimension
               , in_data  # any dimension
               , rigor=False
               , verbose=False)
```

#### 4.1.2 PyTorch wrapper

‘relu()’ PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
relu( input # any dimension
      , rigor=False
      , verbose=False)
```

### 4.2 Leaky ReLU

‘ActivationLeakyReLU()’ applies Rectified Linier Unit over an input data.

```
#include "activation_leakyrelu.hpp"

template<class TYPE =float>
void ActivationLeakyReLU
(
    TYPE      *out_data
    , const TYPE *in_data
    , const uint32_t size
    , const uint16_t channel
    , const uint32_t negative_slope=0x3DCCCCCD
    #if !defined(__SYNTHESIS__)
    , const int rigor=0
    , const int verbose=0
    #endif
);
```

Templates:

- ✧ DTYPE: data type of out\_data, in\_data
  - int, float, double, and so on<sup>3</sup>.

Arguments:

- ✧ out\_data: resultant data in contiguous memory and can be any dimension
  - channel x size x size
- ✧ in\_data: input array in contiguous memory and can be any dimension
  - channel x size x size

<sup>3</sup> Fixed-point data type will be supported.

Future Design Systems	FDS-TD-2020-09-003

- ✧ size: the number of elements per channel
- ✧ channel: the number of channels
- ✧ negative\_slope: apply when in\_data[] is negative; bit pattern of float 0.01 by default.
  - Note that it should carry bit-pattern
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- 'out\_data' and 'in\_array' must have 'channel x size' elements.

#### 4.2.1 Python wrapper

'ActivationLeakyReLu()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
ActivationLeakyReLu( out_data # any dimension
                    , in_data # any dimension
                    , negative_slope
                    , rigor=False
                    , verbose=False)
```

#### 4.2.2 PyTorch wrapper

'leaky\_relu()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
leaky_relu( input # any dimension
           , negative_slope
           , rigor=False
           , verbose=False)
```

### 4.3 Hyperbolic tangent

'ActivationTanh()' applies hyperbolic tangent over an input data.

```
#include "activation_tanh.hpp"

template<class TYPE=float>
void ActivationTanh
(
    TYPE *out_data
    , const TYPE *in_data
    , const uint32_t size
    , const uint16_t channel
    , #if !defined(__SYNTHESIS__)
      const int rigor=0
      , const int verbose=0
    )
```

Future Design Systems	FDS-TD-2020-09-003

```
#endif
);
```

Templates:

- ✧ DTYPE: data type of out\_data, in\_data
  - int, float, double, and so on<sup>4</sup>.

Arguments:

- ✧ out\_data: resultant data in contiguous memory and can be any dimension
  - channel x size x size
- ✧ in\_data: input array in contiguous memory and can be any dimension
  - channel x size x size
- ✧ size: the number of elements per channel
- ✧ channel: the number of channels
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- 'out\_data' and 'in\_array' must have 'channel x size' elements.

#### 4.3.1 Python wrapper

'ActivationTanh()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
ActivationTanh ( out_data # any dimension
                , in_data # any dimension
                , rigor=False
                , verbose=False)
```

#### 4.3.2 PyTorch wrapper

'tanh()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
tanh( input # any dimension
      , rigor=False
      , verbose=False)
```

### 4.4 Sigmoid

'ActivationSigmoid()' applies sigmoid over an input data.

---

<sup>4</sup> Fixed-point data type will be supported.

```
#include "activation_sigmoid.hpp"

template<class TYPE=float>
void ActivationSigmoid
(
    TYPE      *out_data
    , const TYPE  *in_data
    , const uint32_t size
    , const uint16_t channel
    #if !defined(__SYNTHESIS__)
    , const int    rigor=0
    , const int    verbose=0
    #endif
);
```

#### Templates:

- ✧ DTYPE: data type of out\_data, in\_data
  - int, float, double, and so on<sup>5</sup>.

#### Arguments:

- ✧ out\_data: resultant data in contiguous memory and can be any dimension
  - channel x size x size
- ✧ in\_data: input array in contiguous memory and can be any dimension
  - channel x size x size
- ✧ size: the number of elements per channel
- ✧ channel: the number of channel
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- 'out\_data' and 'in\_array' must have 'channel x size' elements.

#### 4.4.1 Python wrapper

'ActivationSigmoid()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
ActivationSigmoid( out_data # any dimension
                  , in_data # any dimension
                  , rigor=False
                  , verbose=False)
```

#### 4.4.2 PyTorch wrapper

'sigmoid()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

---

<sup>5</sup> Fixed-point data type will be supported.

```
sigmoid( input # any dimension
        , rigor=False
        , verbose=False)
```

## 5 Concatenation

### 5.1 Concat 2D

‘Concat2d()’ combines two 2-dimensional arrays.

```
#include "concat_2d.hpp"

template<class TYPE=float>
void Concat2d
(
    TYPE      *out_data
  , const TYPE *in_dataA
  , const TYPE *in_dataB
  , const uint16_t in_rowsA
  , const uint16_t in_colsA
  , const uint16_t in_rowsB
  , const uint16_t in_colsB
  , const uint8_t  dim
  , #if !defined(__SYNTHESIS__)
    const int      rigor=0
    const int      verbose=0
  , #endif
);
```

Templates:

- ✧ DTYPE: data type of out\_data, in\_data
  - int, float, double, and so on<sup>6</sup>.

Arguments:

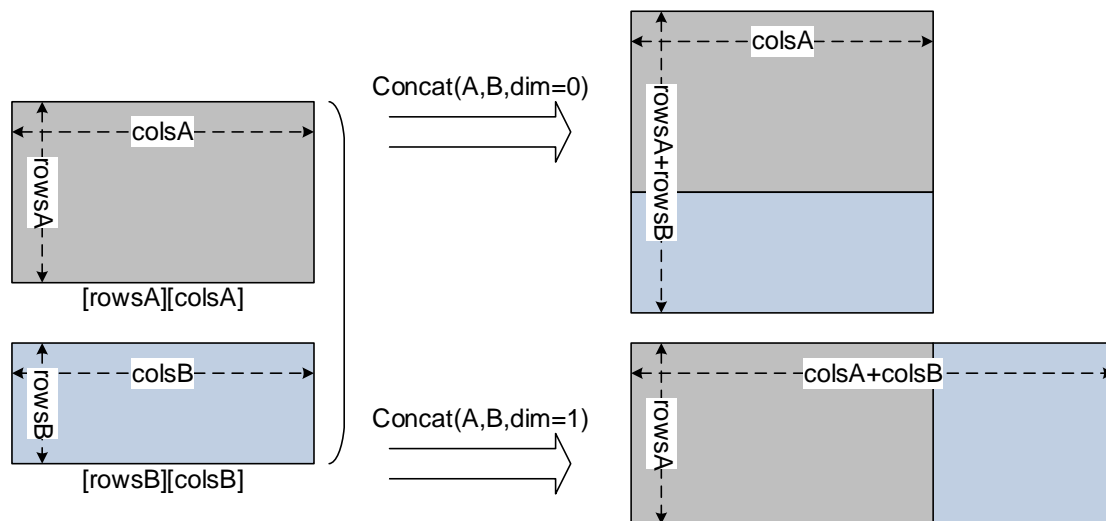
- ✧ out\_data[[]]: resultant 2D array
  - in\_rowsA x (in\_colsA+in\_colsB) when dim is 1
  - (in\_rowsA+in\_rowsB) x in\_colsA when dim is 0
- ✧ in\_dataA[in\_rowA][in\_colA]: input array (rows:height, cols:width)
- ✧ in\_dataB[in\_rowB][in\_colB]: input array (rows:height, cols:width)
- ✧ dim: append array B at the end of array A when dim is 0; append each row of array B at the end of each row of array A when dim is 1.
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- ‘out\_data’, ‘in\_arrayA’ and ‘in\_dataB’ are 2 dimensional array.
- ‘dim’ can be 0 or 1.

---

<sup>6</sup> Fixed-point data type will be supported.



**Figure 4: Concat2d example**

#### 5.1.1 Python wrapper

'Concat2d()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
Concat2d( out_data # in_rowsA x (in_colsA+in_colsB) when dim is 1
          # (in_rowsA+in_rowsB) x in_colsA when dim is 0
          , in_dataA # in_rowsA x in_colsA
          , in_dataB # in_rowsB x in_colsB
          , dim=0    # out_size
          , rigor=False
          , verbose=False)
```

#### 5.1.2 PyTorch wrapper

'cat()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
cat( (inputA
      , inputB
      , dim=0
      , rigor=False
      , verbose=False)
```

## 6 Convolution

### 6.1 Convolution2d

‘Convolution2d()’ applies a 2D convolution over an input data composed of several planes (channels).

```
#include "convolution_2d.hpp"

template<class TYPE=float>
void Convolution2d
(
    TYPE      *out_data
  , const TYPE *in_data
  , const TYPE *kernel
  , const TYPE *bias
  , const uint16_t out_size
  , const uint16_t in_size
  , const uint8_t  kernel_size
  , const uint16_t bias_size
  , const uint16_t in_channel
  , const uint16_t out_channel
  , const uint8_t  stride
  , const uint8_t  padding=0
  #if !defined(__SYNTHESIS__)
  , const int      rigor=0
  , const int      verbose=0
  #endif
)

```

Templates:

- ✧ DTYPE: data type of out\_data, in\_data, kernel, bias
  - int, float, double, and so on<sup>7</sup>.

Arguments:

- ✧ out\_data[out\_channel][out\_size][out\_size] resultant array in square
- ✧ in\_data[in\_channel][in\_size][in\_size] input feature array in square
- ✧ kernel[in\_channel][out\_channel][kernel\_size][kernel\_size] filter array in square
- ✧ bias[bias\_size]: bias for output and each element corresponds each output channel
- ✧ out\_size: the number of elements in width direction (column)
- ✧ in\_size: the number of elements in width direction (column)
- ✧ kernel\_size: the number of elements in width direction (column)
- ✧ bias\_size: the number of elements in bias (zero by default)
- ✧ in\_channel: the number of input channels
- ✧ out\_channel: the number of output channels, i.e., filters; the number of channels produced by the convolution.
- ✧ strid: stride of the convolution (one by default)
- ✧ padding: the number of zero-paddings on both sides (zero by default, not implemented yet)
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

---

<sup>7</sup> Fixed-point data type will be supported.

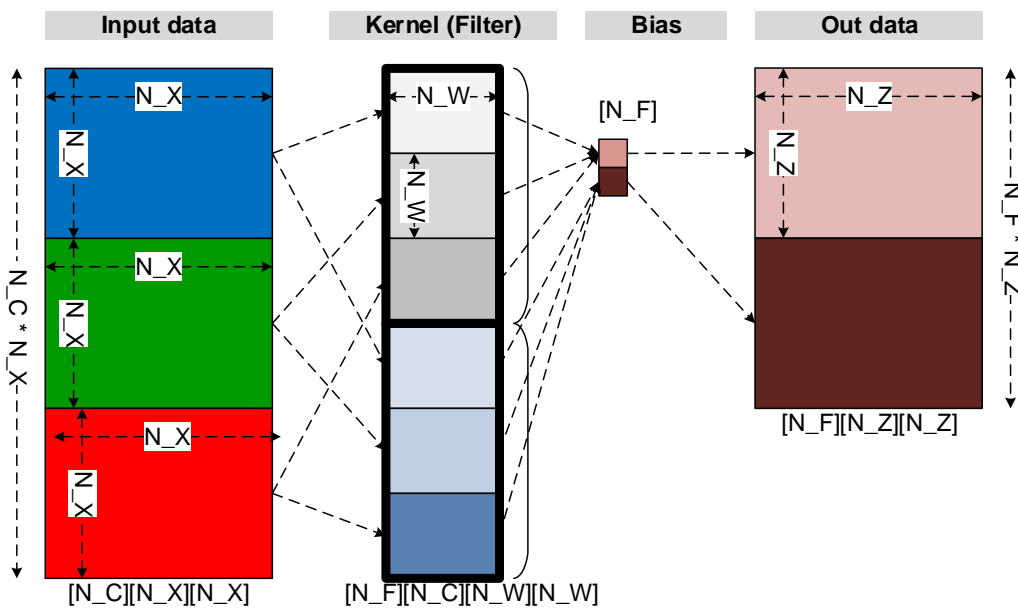
'minibatch', 'dilation', and 'goups' are not supported, but 'minibatch' is supported through PyTorch wrapper.

Array size must keep following rules.

- 'out\_data', 'in\_data', 'kernel' are square matrix
- 'out\_size' =  $((in\_size - kernel\_size + 2 * padding) / stride) + 1$
- 'kernel\_size' should be odd number
- 'bias\_size' should be the same as 'out\_channel', i.e., the number of kernels
- 'stride' should be  $\geq 1$
- 'padding' should be  $\geq 0$  and will be  $\text{floor}(kernel\_size / 2)$  normal case

Each data arrays are row-major contiguous and kernel array can be seen as 4-dimensional array of  $[out\_channel \times in\_channel \times kernel\_size \times kernel\_size]$ .

Figure 5 shows an example case of three input channels and two output channels<sup>8</sup>, in which biases are two and kernels are 'out\_channel x in\_channel'.



**Figure 5: Convolution2d example with 3 input channels, 2 output channels**

### 6.1.1 Python wrapper

'Convolution2d()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.



Future Design Systems	FDS-TD-2020-09-003

```
Convolution2d( out_data # out_channel x out_size x out_size
               , in_data  # in_channel x in_size x in_size
               , kernel   # out_channel x in_channel x kernel_size x kernel_size
               , stride=1
               , padding=0
               , bias=0
               , rigor=False
               , verbose=False)
```

### 6.1.2 PyTorch wrapper

'conv2d()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
conv2d( input # minibatch x in_channel x in_size x in_size
        , weight # out_channel x in_channel x kernel_size x kernel_size
        , bias=None # out_channel
        , stride=1
        , padding=0
        , dilation=1
        , groups=1
        , rigor=False
        , verbose=False)
```

## 7 Deconvolution (Transposed convolution)

### 7.1 deconvolution2d

'DeConvolution2d()' applies a 2D transposed convolution over an input signal composed of several input planes (channels).

```
#include "deconvolution_2d.hpp"
```

```
template<class TYPE=float>
void Deconvolution2d
(
    TYPE      *out_data
  , const TYPE *in_data
  , const TYPE *kernel
  , const TYPE *bias
  , const uint16_t out_size
  , const uint16_t in_size
  , const uint8_t  kernel_size
  , const uint8_t  bias_size
  , const uint16_t in_channel
  , const uint16_t out_channel
  , const uint8_t  stride
  , const uint8_t  padding=0
  , const int      rigor=0
  , const int      verbose=0
  , const int      __SYNTHESIS__
)
#endif
```

Future Design Systems	FDS-TD-2020-09-003

```
);
```

Templates:

- ✧ DTYPE: data type of in\_data, kernel, bias, out\_data; int, float, double, and so on<sup>9</sup>.

Arguments:

- ✧ out\_data[out\_channel][out\_size][out\_size] resultant array in square
- ✧ in\_data[in\_channel][in\_size][in\_size] input feature array in square
- ✧ kernel[in\_channel][out\_channel][kernel\_size][kernel\_size] filter array in square
- ✧ out\_size: the number of elements in width direction (column)
- ✧ in\_size: the number of elements in width direction (column)
- ✧ kernel\_size: the number of elements in width direction (column)
- ✧ in\_channel: the number of input channels
- ✧ out\_channel: the number of output channels, i.e., filters; the number of channels produced by the convolution.
- ✧ stride: stride of the convolution (one by default)
- ✧ padding: the number of zero-paddings on both sides (zero by default, not implemented yet)
- ✧ bias[bias\_size]: bias for output and each element corresponds each output channel
- ✧ bias\_size: the number of elements in bias (zero by default)

'padding', 'minibatch', 'dilation', 'goups' are not supported yet.

Array size must keep following rules.

- 'out\_data', 'in\_data', 'kernel' are square matrix
- 'out\_size' = (((in\_size-kernel\_size+2\*padding)/stride)+1)
- 'kernel\_size' should be odd number
- 'stride' should be positive and larger than 0 (should not be 0)
- 'padding' should be positive including 0 and will be floor(kernel\_size/2)
- 'bias\_size' should be 'out\_channel'.

### 7.1.1 Python wrapper

'Deconvolution2d()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
Deconvolution2d( out_data # any dimension
                 , in_data # any dimension
                 , weight
                 , bias=None
                 , stride=1
                 , padding=0
                 , output_padding=0
```

<sup>9</sup> Fixed-point data type will be supported.

Future Design Systems	FDS-TD-2020-09-003

```
, groups=1
, dilation=1
, rigor=False
, verbose=False)
```

### 7.1.2 PyTorch wrapper

'conv\_transpose2d()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
conv_transpose2d( input # any dimension
, weight
, bias=None
, stride=1
, padding=0
, output_padding=0
, groups=1
, dilation=1
, rigor=False
, verbose=False)
```

## 8 Linear (Fully connected)

### 8.1 Linear 1D

'Linear1d()' applies a 1D vector multiplication over an input data.

```
#include "linear_1d.hpp"

template< class TYPE=float
, int ReLu=0
, const int LeakyReLu=0
, const uint32_t negative_slope=0x3DCCCCCD
>
void Linear1d
(
    TYPE *out_data
, const TYPE *in_data
, const TYPE *weight
, const TYPE *bias
, const uint16_t out_size
, const uint16_t in_size
, const uint16_t bias_size
, #if !defined(__SYNTHESIS__)
, const int rigor=0
, const int verbose=0
, #endif
);
```

Templates:

✧ DTYPE: data type of out\_data, in\_data, kernel, bias, out\_data

Future Design Systems	FDS-TD-2020-09-003

- int, float, double, and so on<sup>10</sup>.
- ✧ ReLu: add ReLU activation at the end of linear operation (it should be 0 when LeakyReLU is 1)
- ✧ LeakyReLu: Leaky-ReLU activation at the end of linear operation (it should be 0 when ReLU is 1)
- ✧ negative\_slope: Slope to be applied for negative input with 0.001 by default and only valid when LeakyReLU is 1.

Arguments:

- ✧ out\_data[out\_size]: resultant 1D array
- ✧ in\_data[in\_size]: input feature 1D array
- ✧ weight[out\_size][in\_size]: weights
  - weight should be transposed before multiplying to in\_data
- ✧ bias[out\_size]: bias 1D array
- ✧ out\_size: the number of elements of out\_data[]
- ✧ in\_size: the number of elements of in\_data[]
- ✧ bias\_size: the number of elements of bias[]
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- 'out\_data' and 'in\_data' are 1 dimensional array.
- 'weight' is 2 dimensional array
- 'weight' should be 'out\_size x in\_size'
- 'bias' should be 1 dimensional array
- 'bias\_size' should be the same as 'out\_size'

### 8.1.1 Python wrapper

'Linear1d()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
Linear1d( out_data  # out_size
         , in_data  # in_size
         , weight # out_size x in_size
         , bias # out_size
         , rigor=False
         , verbose=False)
```

### 8.1.2 PyTorch wrapper

'linear()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

---

<sup>10</sup> Fixed-point data type will be supported.

```
linear( input # minibatch x 1 x in_size
        , weight # out_size x in_size
        , bias # out_size
        , rigor=False
        , verbose=False)
```

## 8.2 Linear ND

‘LinearNd()’ applies a N-D matrix multiplication over an input data.

```
#include "linear_nd.hpp"

template<class TYPE=float, int ReLu=0>
void LinearNd
(
    TYPE *out_data
    , const TYPE *in_data
    , const TYPE *weight
    , const TYPE *bias
    , const uint16_t out_size
    , const uint16_t in_size
    , const uint16_t bias_size
    , const uint8_t ndim
    #if !defined(__SYNTHESIS__)
    , const int rigor=0
    , const int verbose=0
    #endif
);
```

Templates:

- ✧ DTYPE: data type of out\_data, in\_data, kernel, bias, out\_data
  - int, float, double, and so on<sup>11</sup>.
- ✧ ReLu: add ReLU activation at the end of linear operation

Arguments:

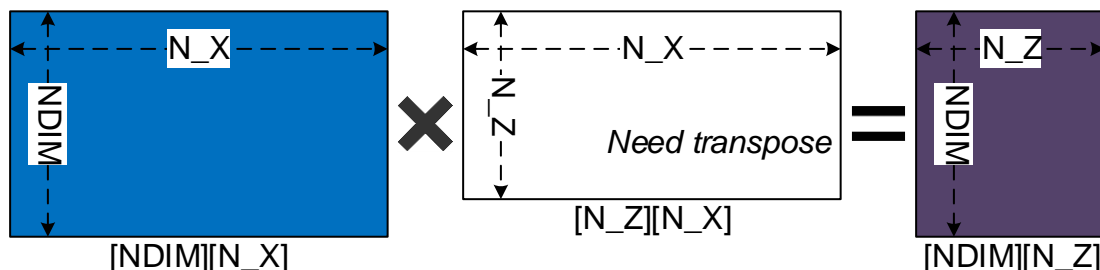
- ✧ out\_data[ndim][out\_size]: resultant 1D array
- ✧ in\_data[ndim][in\_size]: input feature 1D array
- ✧ weight[out\_size][in\_size]: weight
  - weight should be transposed before multiplying to in\_data
- ✧ bias[out\_size]: bias 1D array
- ✧ out\_size: the number of elements of out\_data[]
- ✧ in\_size: the number of elements of in\_data[]
- ✧ bias\_size: the number of elements of bias[]
- ✧ ndim: dimension of in\_data
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

Array size must keep following rules.

- ‘out\_data’ and ‘in\_data’ are 2 dimensional array.

<sup>11</sup> Fixed-point data type will be supported.

- 'weight' is 2 dimensional array
- 'weight' should be 'out\_size x in\_size'
- 'bias' should be 1 dimensional array
- 'bias\_size' should be the same as 'out\_size'



**Figure 6: Matrix multiplication over 2 dimensional array case**

### 8.2.1 Python wrapper

'LinearNd()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
LinearNd( out_data # ndim x out_size
         , in_data  # ndim x in_size
         , weight   # out_size x in_size
         , bias     # out_size
         , rigor=False
         , verbose=False)
```

### 8.2.2 PyTorch wrapper

'linearNd()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
linearNd( input  # minibatch x ndim x in_size
         , weight # out_size x in_size
         , bias   # out_size
         , rigor=False
         , verbose=False)
```

## 9 Normalization

### 9.1 Batch normalization 2D

'Batchnorm()' applies a batch normalization for each channel across a batch of data.

$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

where 'y' is outut, 'x' is input, 'E[x]' is mean, 'Var[x]' is variance, 'gamma' is scaling factor, 'beta' is shift factor (bias), 'epsilon' is value for numerical stability.

```
#include "norm_2d_batch.hpp"

template< class TYPE=float
        , int LeakyReLu=0
        , int negative_slope1000=100>
void Norm2dBatch
(
    TYPE      *out_data
  , const TYPE *in_data
  , const TYPE *running_mean
  , const TYPE *running_var
  , const TYPE *scale
  , const TYPE *bias
  , const uint32_t in_size
  , const uint16_t scale_size
  , const uint16_t bias_size
  , const uint16_t in_channel
  , const float   epsilon=1E-5
  #if !defined(__SYNTHESIS__)
  , const int     rigor=0
  , const int     verbose=0
  #endif
);
```

Templates:

- ✧ DTYPE: data type of out\_data, in\_data, running\_mean, running var, scale and bias
  - int, float, double, and so on.
- ✧ LeakyReLu: apply Leaky ReLU at the end of this operation
- ✧ negative\_slope1000: 100 means 0.1

Arguments:

- ✧ out\_data[in\_channel][in\_size][in\_size] resultant array in square
- ✧ in\_data[in\_channel][in\_size][in\_size] input feature array in square
- ✧ running\_mean: running\_mean[in\_channel], mean of each channel
- ✧ running\_var: running\_var[in\_channel], standard variance of each channel
- ✧ scale: scale[in\_channel], scale factor (by default, 1)
- ✧ bias: bias[in\_channel], shift factor (by default, 0)
- ✧ in\_size: the number of elements in width direction (column)
- ✧ scale\_size: the number of elements in scale[], it can be 0 or in\_channel
- ✧ bais\_sie: the number of elements in bias[], it can be 0 or in\_channel
- ✧ in\_channel: the number of input channels
- ✧ epsilon: numeric stability (by default, 1e-5)

Future Design Systems	FDS-TD-2020-09-003

Array size must keep following rules.

- 'out\_data' and 'in\_data' are the same size
- scale\_size, bias\_size can be 0 or in\_channel.

### 9.1.1 Python wrapper

'Norm2dBatch()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
Norm2dBatch( out_data  # out_size
            , in_data   # in_size
            , rigor=False
            , verbose=False)
```

### 9.1.2 PyTorch wrapper

'batch\_norm()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
batch_norm( input # minibatch x 1 x in_size
           , rigor=False
           , verbose=False)
```

## 10 Pooling

### 10.1 Pooling2dMax

'Pooling2dMax()' applies a 2D max pooling over an input data composed of several input planes (channels).

```
#include "pooling_2d_max.hpp"

template< class TYPE=float
        , const int ReLu=0
        , const int LeakyReLu=0
        , const uint32_t negative_slope=0x3DCCCCCD
        >
void Pooling2dMax
(
    TYPE *out_data
    , const TYPE *in_data
    , const uint16_t out_size
    , const uint16_t in_size
    , const uint8_t kernel_size
    , const uint16_t channel
    , const uint8_t stride
    , const uint8_t padding=0
    , const int ceil_mode=0
```



```

    #if !defined(__SYNTHESIS__)
    , const int    rigor=0
    , const int    verbose=0
    #endif
);

```

#### Templates:

- ✧ DTYPE: data type of out\_data, in\_data, kernel, bias, out\_data
  - int, float, double, and so on<sup>12</sup>.
- ✧ ReLU: apply ReLU at the end of this operation when 1
- ✧ LeakyReLU: apply leaky ReLU at the end of this operation when 1
- ✧ negative\_slope: Slope to be applied for negative input and only valid when LeakyReLU is 1. 0.01 by default.

#### Arguments:

- ✧ out\_data[channel][out\_size][out\_size] resultant array in square
- ✧ in\_data[channel][in\_size][in\_size] input feature array in square
- ✧ out\_size: the number of elements in width direction (column)
- ✧ in\_size: the number of elements in width direction (column)
- ✧ kernel\_size: the number of elements in width direction (column)
- ✧ channel: the number of input/output channels
- ✧ stride: stride of the convolution (one by default)
- ✧ padding: the number of zero-paddings on both sides (zero by default, not implemented yet)
- ✧ ceil\_mode: use ceil() instead of floor to calculate out\_size
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

‘minibatch’ is not supported.

Array size must keep following rules.

- ‘out\_data’, ‘in\_data’, ‘kernel’ are square matrix
- ‘out\_size’ = floor(((in\_size-kernel\_size+2\*padding)/stride)+1) when ‘ceil\_mode’ is 0 by default. Otherwise, ceil() is used instead.
- ‘kernel\_size’ should be odd number
- ‘stride’ should be >=1
- ‘padding’ should be >=0 and will be floor(kernel\_size/2) normal case

#### 10.1.1 Python wrapper

‘Pooling2dMax()’ wrapper gets NumPy arguments for array and ‘out\_data’ carries calculated result. It returns ‘True’ on success or ‘False’ on failure.

<sup>12</sup> Fixed-point data type will be supported.

```
Pooling2dMax( out_data # out_channel x out_size x out_size
              , in_data  # in_channel x in_size x in_size
              , kernel_size
              , stride=1
              , padding=0
              , ceil_mode=False
              , rigor=False
              , verbose=False)
```

### 10.1.2 PyTorch wrapper

'max\_pool2d()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

```
max_pool2d( input # minibatch x in_channel x in_size x in_size
            , kernel_size
            , stride=1
            , padding=0
            , ceil_mode=False
            , groups=1
            , rigor=False
            , verbose=False)
```

## 10.2 Pooling2dAvg

'Pooling2dAvg()' applies a 2D average pooling over an input data composed of several input planes (channels).

```
#include "pooling_2d_avg.hpp"

template<class TYPE=float>
void Pooling2dAvg
(
    TYPE *out_data
    , const TYPE *in_data
    , const uint16_t out_size
    , const uint16_t in_size
    , const uint8_t kernel_size
    , const uint16_t channel
    , const uint8_t stride
    , const uint8_t padding=0
    , const int ceil_mode=0
    , #if !defined(__SYNTHESIS__)
      const int rigor=0
      , const int verbose=0
    #endif
)
)
```

Templates:

✧ DTYPE: data type of out\_data, in\_data, kernel, bias, out\_data

Future Design Systems	FDS-TD-2020-09-003

□ int, float, double, and so on<sup>13</sup>.

Arguments:

- ✧ out\_data[channel][out\_size][out\_size] resultant array in square
- ✧ in\_data[channel][in\_size][in\_size] input feature array in square
- ✧ out\_size: the number of elements in width direction (column)
- ✧ in\_size: the number of elements in width direction (column)
- ✧ kernel\_size: the number of elements in width direction (column)
- ✧ channel: the number of input/output channels
- ✧ stride: stride of the convolution (one by default)
- ✧ padding: the number of zero-paddings on both sides (zero by default, not implemented yet)
- ✧ ceil\_mode: use ceil() instead of floor to calculate out\_size
- ✧ rigor: check values rigorously when 1
- ✧ verbose: print more message when 1

'minibatch' is not supported.

Array size must keep following rules.

- 'out\_data', 'in\_data', 'kernel' are square matrix
- 'out\_size' = floor(((in\_size-kernel\_size+2\*padding)/stride)+1) when 'ceil\_mode' is 0 by default. Otherwise, ceil() is used instead.
- 'kernel\_size' should be odd number
- 'stride' should be >=1
- 'padding' should be >=0 and will be floor(kernel\_size/2) normal case

### 10.2.1 Python wrapper

'Pooling2dAvg()' wrapper gets NumPy arguments for array and 'out\_data' carries calculated result. It returns 'True' on success or 'False' on failure.

```
Pooling2dAvg( out_data  # out_channel x out_size x out_size
              , in_data  # in_channel x in_size x in_size
              , kernel_size
              , stride=1
              , padding=0
              , ceil_mode=False
              , rigor=False
              , verbose=False)
```

### 10.2.2 PyTorch wrapper

'avg\_pool2d()' PyTorch wrapper gets PyTorch tensor arguments for array and returns calculated result. It calls Python wrapper after converting PyTorch tensor to NumPy array.

<sup>13</sup> Fixed-point data type will be supported.

```
avg_pool2d( input # minibatch x in_channel x in_size x in_size
            , kernel_size
            , stride=1
            , padding=0
            , ceil_mode=False
            , groups=1
            , rigor=False
            , verbose=False)
```

## 11 Project: LeNet-5

### 11.1 LeNet-5

LeNet-5 is a convolutional network designed for handwritten and machine-printed character recognition. Figure 7 shows network and its data size. Table 1 shows a summary of the network and its parameters.

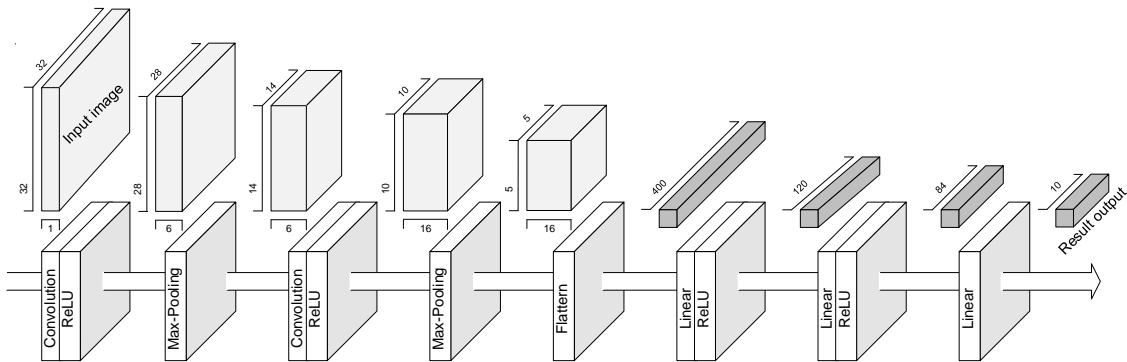


Figure 7: LeNet-5 network

Table 1: LeNet-5 summary

Layer		feature map (channel x W x H)	kernel size	stride	padding	Activation	Parameters (bias+weight)	Feature map
input	image	1 32x32	.	.	.	.	.	1,024
1	CONV	6 28x28	5x5	1	0	ReLU	6+150	4,704
2	MaxPOOL	6 14x14	2x2	2	0	.	0	1,176
3	CONV	16 10x10	5x5	1	0	ReLU	16+2,400	1,600
4	MaxPOOL	16 5x5	2x2	2	0	.	0	400
5	FC	120 1x1	5x5	1	0	ReLU	120+48,000	120
6	FC	. 84	.	.	.	ReLU	84+10,080	84
7	FC	. 10	.	.	.	softmax	0	10
							61,706	9,118

#### 11.1.1 LeNet-5 C++ version

Following code shows an example implementation of LeNet-5 network model using DLR deep learning routines, where 'lenet5\_params.h' should contain

trained results of weights and biases. Flatten step before the first linear layer is not required for DLR since its data array can be seen linear array.

```

#ifndef DTYPE
#define DTYPE float
#endif

#include "dlr.hpp"
#include "lenet5_params.h" // come from LeNet-5.pytorch

void lenet5(    DTYPE classes[10]
              , const DTYPE image [32][32]
              #if !defined(__SYNTHESIS__)
              , const int rigor
              , const int verbose
              #endif
              )
{
    DTYPE    c1_out_data[6][28][28];
    const DTYPE (*c1_in_data)[32][32]=(DTYPE (*)[32][32])image; // [1][32][32]
    const DTYPE (*c1_kernel)[1][5][5]=(DTYPE (*)[1][5][5])conv1_weight; // [6][1][5][5]
    const DTYPE (*c1_bias)=conv1_bias; // [6]
    const uint16_t c1_out_size=28;
    const uint16_t c1_in_size=32;
    const uint8_t c1_kernel_size=5;
    const uint16_t c1_bias_size=6;
    const uint16_t c1_in_channel=1;
    const uint16_t c1_out_channel=6;
    const uint8_t c1_stride=1;
    const uint8_t c1_padding=0;

    Convolution2d<DTYPE> (
        (DTYPE *)c1_out_data
        , (DTYPE *)c1_in_data
        , (DTYPE *)c1_kernel
        , (DTYPE *)c1_bias
        , c1_out_size
        , c1_in_size
        , c1_kernel_size
        , c1_bias_size
        , c1_in_channel
        , c1_out_channel
        , c1_stride
        , c1_padding
        #if !defined(__SYNTHESIS__)
        , rigor
        , verbose
        #endif
    );

    DTYPE    p1_out_data[6][14][14];
    const DTYPE (*p1_in_data)[28][28]=c1_out_data; // [6][28][28]
    const uint16_t p1_out_size=14;
    const uint16_t p1_in_size=28;
    const uint8_t p1_kernel_size=2;
    const uint8_t p1_channel=6;
    const uint8_t p1_stride=2;
    const uint8_t p1_padding=0;
    const int p1_ceil_mode=0;

    Pooling2dMax<DTYPE, 1> ( // ReLU embedded
        (DTYPE *)p1_out_data // out_channel x out_size x out_size
        , (DTYPE *)p1_in_data // in_channel x in_size x in_size

```

```

,      p1_out_size // only for square matrix
,      p1_in_size  // only for square matrix
,      p1_kernel_size // only for square matrix
,      p1_channel
,      p1_stride
,      p1_padding
,      p1_ceil_mode
#if !defined(__SYNTHESIS__)
,      rigor
,      verbose
#endif
);

DTYPE   c2_out_data[16][10][10];
const DTYPE (*c2_in_data)[14][14]=p1_out_data; // [6][14][14]
const DTYPE (*c2_kernel)[6][5][5]=(DTYPE *) [6][5][5]conv2_weight; // [16][6][5][5]
const DTYPE (*c2_bias)=conv2_bias; // [16]
const uint16_t c2_out_size=10;
const uint16_t c2_in_size=14;
const uint8_t  c2_kernel_size=5;
const uint16_t c2_bias_size=16;
const uint16_t c2_in_channel=6;
const uint16_t c2_out_channel=16;
const uint8_t  c2_stride=1;
const uint8_t  c2_padding=0;

Convolution2d<DTYPE> (
    (DTYPE *)c2_out_data
, (DTYPE *)c2_in_data
, (DTYPE *)c2_kernel
, (DTYPE *)c2_bias
,      c2_out_size
,      c2_in_size
,      c2_kernel_size
,      c2_bias_size
,      c2_in_channel
,      c2_out_channel
,      c2_stride
,      c2_padding
#if !defined(__SYNTHESIS__)
,      rigor
,      verbose
#endif
);

DTYPE   p2_out_data[16][5][5];
const DTYPE (*p2_in_data)[10][10]=c2_out_data; // [16][10][10]
const uint16_t p2_out_size=5;
const uint16_t p2_in_size=10;
const uint8_t  p2_kernel_size=2;
const uint8_t  p2_channel=16;
const uint8_t  p2_stride=2;
const uint8_t  p2_padding=0;
const int      p2_ceil_mode=0;

Pooling2dMax<DTYPE, 1> ( // ReLU embedded
    (DTYPE *)p2_out_data
, (DTYPE *)p2_in_data
,      p2_out_size
,      p2_in_size
,      p2_kernel_size
,      p2_channel
,      p2_stride
,      p2_padding

```

```

,      p2_ceil_mode
#if !defined(__SYNTHESIS__)
,      rigor
,      verbose
#endif
);

    DTYPE  f1_out_data[120];
const DTYPE (*f1_in_data)[5][5]=p2_out_data; // [16][5][5]
const DTYPE (*f1_weight)[400]=(DTYPE *) [400]fc1_weight; // [120][400]
const DTYPE (*f1_bias)=fc1_bias; // [120]
const uint16_t f1_out_size=120;
const uint16_t f1_in_size=16*5*5;
const uint16_t f1_bias_size=120;

Linear1d<DTYPE, 1> ( // ReLU embedded
    (DTYPE *)f1_out_data
    , (DTYPE *)f1_in_data
    , (DTYPE *)f1_weight
    , (DTYPE *)f1_bias
    ,      f1_out_size
    ,      f1_in_size
    ,      f1_bias_size
    #if !defined(__SYNTHESIS__)
    ,      rigor
    ,      verbose
    #endif
);

    DTYPE  f2_out_data[84];
const DTYPE (*f2_in_data)=f1_out_data; // [120]
const DTYPE (*f2_weight)[120]=(DTYPE *) [120]fc2_weight; // [84][120]
const DTYPE (*f2_bias)=fc2_bias; // [84]
const uint16_t f2_out_size=84;
const uint16_t f2_in_size=120;
const uint16_t f2_bias_size=84;

Linear1d<DTYPE, 1> ( // ReLU embedded
    (DTYPE *)f2_out_data
    , (DTYPE *)f2_in_data
    , (DTYPE *)f2_weight
    , (DTYPE *)f2_bias
    ,      f2_out_size
    ,      f2_in_size
    ,      f2_bias_size
    #if !defined(__SYNTHESIS__)
    ,      rigor
    ,      verbose
    #endif
);

    DTYPE (*f3_out_data)=classes; // [10]
const DTYPE (*f3_in_data)=f2_out_data; // [84]
const DTYPE (*f3_weight)[84]=(DTYPE *) [84]fc3_weight; // [10][84]
const DTYPE (*f3_bias)=fc3_bias; // [10]
const uint16_t f3_out_size=10;
const uint16_t f3_in_size=84;
const uint16_t f3_bias_size=10;

Linear1d<DTYPE, 0> ( // ReLU not embedded
    (DTYPE *)f3_out_data
    , (DTYPE *)f3_in_data
    , (DTYPE *)f3_weight
    , (DTYPE *)f3_bias

```

```

        ,      f3_out_size
        ,      f3_in_size
        ,      f3_bias_size
    #if !defined(__SYNTHESIS__)
        ,      rigor
        ,      verbose
    #endif
    );
}

```

### 11.1.2 LeNet-5 C PyTorch version

Following code shows an example of LeNet-5 network using DLR PyTorch wrapper functions.

```

import torch
import python.torch as dlr

def lenet5_infer( input
    , c1_kernel, c1_bias
    , c2_kernel, c2_bias
    , f1_weight, f1_bias
    , f2_weight, f2_bias
    , f3_weight, f3_bias
    , softmax=True
    , pkg=dlr
    , rigor=False):
    """ LeNet-5 inference network """
    z = pkg.conv2d( input=input, weight=c1_kernel, bias=c1_bias, stride=1
        , padding=0, dilation=1, groups=1, rigor=rigor )
    z = pkg.relu( input=z, rigor=rigor )
    z = pkg.max_pool2d( input=z, kernel_size=2, stride=2, padding=0,
        ceil_mode=False, rigor=rigor )

    z = pkg.conv2d( input=z, weight=c2_kernel, bias=c2_bias, stride=1
        , padding=0, dilation=1, groups=1, rigor=rigor )
    z = pkg.relu( input=z, rigor=rigor )
    z = pkg.max_pool2d( input=z, kernel_size=2, stride=2, padding=0,
        ceil_mode=False, rigor=rigor )

    z = z.view(z.shape[0], -1)
    z = pkg.linear( input=z, weight=f1_weight, bias=f1_bias, rigor=rigor )
    z = pkg.relu( input=z, rigor=rigor )
    z = pkg.linear( input=z, weight=f2_weight, bias=f2_bias, rigor=rigor )
    z = pkg.relu( input=z, rigor=rigor )
    z = pkg.linear( input=z, weight=f3_weight, bias=f3_bias, rigor=rigor )
    if softmax: z = F.softmax( input=z, dim=1)
    else: z = pkg.relu( input=z )
    return z

```

## 12 Troubleshooting

To be added.



Future Design Systems	FDS-TD-2020-09-003

## 13 Acknowledgement

This work was supported partially by KARI<sup>14</sup> (Korea Aerospace Research Institute) under the “Development of deep learning hardware accelerator for microsatellite<sup>15</sup>” (Contract 2019110A35C-00).

Some part of this work was carried out partially by Handong Global University<sup>16</sup> 2020 Summer Internship program.

Some part of this work was supported partially by ETRI<sup>17</sup> (Electronics and Telecommunications Research Institute) under the “Development and training of FPGA-based AI semiconductor design environment<sup>18</sup>” (Contract EA20202206).

## 14 References

- [1] Future Design Systems, CON-FMC User Manual, FDS-TD-2018-03-001, 2018.
- [2] NumPy, <https://numpy.org>
- [3] Ctypes., <https://docs.python.org/3/library/ctypes.html>
- [4] PyTorch, <https://pytorch.org>

## Wish list

- ☐ Softmax, upsampling

## Index

C		<b>DPU</b>	
	CON-FMC 4		Deep Learning Processing Unit 4
D		H	
	<b>DLIP</b>		<b>HDL</b>
	Deep Learning Intellectual Property 4		Hardware Description Language 4
	<b>DLR</b>		
	Deep Learning Routines 4		

<sup>14</sup> 한국항공우주연구원, <https://www.kari.re.kr>

<sup>15</sup> 초소형위성용 딥러닝 알고리즘 하드웨어 가속기 개발

<sup>16</sup> 한동대학교, <https://www.handong.edu>, Chaeon Lim, Yoonseong Lim, Geunsu Song.

<sup>17</sup> 전자통신연구원, <https://www.etri.re.kr>

<sup>18</sup> FPGA 기반 AI 반도체 설계환경 개발 및 교육

Future Design Systems	FDS-TD-2020-09-003

## Revision history

- ☐ 2020.12.17: Minor correction
- ☐ 2020.11.12: License scheme changed from proprietary to open.
- ☐ 2020.10.20: 'channel' argument added for activation functions
- ☐ 2020.10.01: Major update.
- ☐ 2020.09.01: First version released.
- ☐ 2020.08.31: More IP have been added.
- ☐ 2020.06.23: Initial version released
- ☐ 2020.03.10: Started by Ando Ki (adki@future-ds.com)

– End of document –