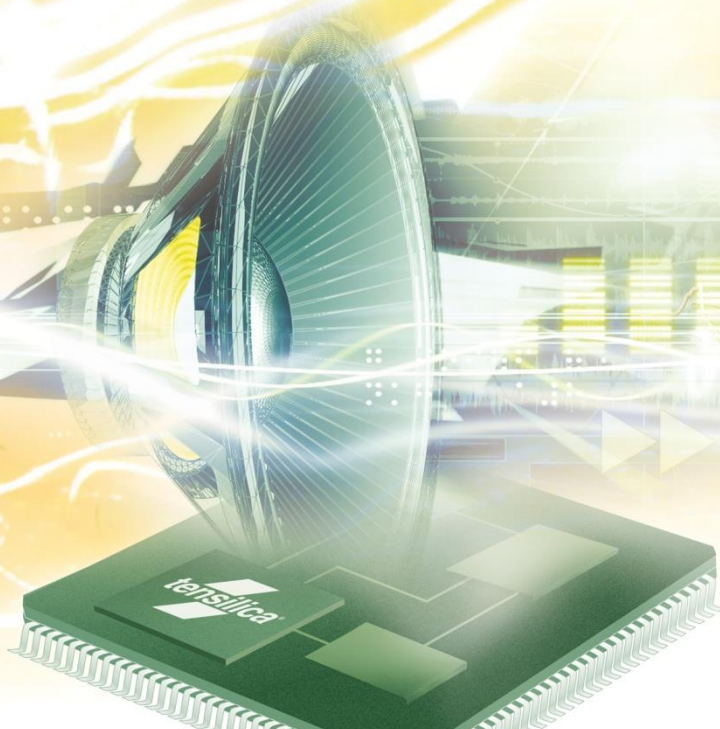




HiFi 5 Neural Network Library

Programmer's Guide

For HiFi DSPs



Cadence Design Systems, Inc.
2655 Seely Ave.
San Jose, CA 95134
www.cadence.com

© 2020 Cadence Design Systems, Inc.
All rights reserved worldwide

This publication is provided “AS IS.” Cadence Design Systems, Inc. (hereafter “Cadence”) does not make any warranty of any kind, either expressed or implied, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose. Information in this document is provided solely to enable system and software developers to use our processors. Unless specifically set forth herein, there are no express or implied patent, copyright or any other intellectual property rights or licenses granted hereunder to design or fabricate Cadence integrated circuits or integrated circuits based on the information in this document. Cadence does not warrant that the contents of this publication, whether individually or as one or more groups, meets your requirements or that the publication is error-free. This publication could include technical inaccuracies or typographical errors. Changes may be made to the information herein, and these changes may be incorporated in new editions of this publication.

© 2019 Cadence, the Cadence logo, Allegro, Assura, Broadband Spice, CDNLIVE!, Celtic, Chipestimate.com, Conformal, Connections, Denali, Diva, Dracula, Encounter, Flashpoint, FLIX, First Encounter, Incisive, Incyte, InstallScape, NanoRoute, NC-Verilog, OrCAD, OSKit, Palladium, PowerForward, PowerSI, PSpice, Purespec, Puresuite, Quickcycles, SignalStorm, Sigrity, SKILL, SoC Encounter, SourceLink, Spectre, Specman, Specman-Elite, SpeedBridge, Stars & Strikes, Tensilica, TripleCheck, TurboXim, Vectra, Virtuoso, VoltageStorm, Xcelium, Xplorer, Xtensa, and Xtreme are either trademarks or registered trademarks of Cadence Design Systems, Inc. in the United States and/or other jurisdictions.

OSCI, SystemC, Open SystemC, Open SystemC Initiative, and SystemC Initiative are registered trademarks of Open SystemC Initiative, Inc. in the United States and other countries and are used with permission. All other trademarks are the property of their respective holders.

Version 1.5
February 2021

Contents

1.	Introduction to the HiFi 5 NN Library	1
1.1	Organization of the HiFi 5 NN Library Package.....	1
1.1.1	Document Overview.....	2
1.2	HiFi 5 NN Library Specification	2
1.2.1	Low Level Kernels.....	2
1.2.2	Layers	3
1.2.3	Support for TensorFlow Lite Micro Operators.....	3
2.	Generic HiFi NN Layer API	5
2.1	Shape.....	5
2.2	Memory Management.....	7
2.2.1	API Handle / Persistent Memory.....	7
2.2.2	Scratch Memory	7
2.2.3	Weights and Biases Memory	7
2.2.4	Input Buffer	7
2.2.5	Output Buffer.....	7
2.3	Generic API Errors	8
2.3.1	Common API Errors.....	8
2.4	C Language API.....	9
2.4.1	Startup Functions.....	10
2.4.2	Query Functions.....	11
2.4.3	Initialization Functions.....	11
2.4.4	Execution Functions.....	11
3.	HiFi 5 NN Library – Low-Level Kernels	12
3.1	Matrix X Vector Multiplication Kernels	12
3.1.1	Matrix X Vector Kernels	12
3.1.2	Fused (Activation) Matrix X Vector Kernels	16
3.1.3	Matrix X Vector Batch Kernels	20
3.1.4	Matrix Multiplication Kernels	22
3.1.5	Matrix X Vector Kernels with Output Stride	25
3.2	Convolution Kernels.....	27
3.2.1	Standard 2D Convolution Kernel	27
3.2.2	Standard 2D Convolution Kernel with Dilation.....	32
3.2.3	Standard 1D Convolution Kernel	36
3.2.4	Depthwise Separable 2D Convolution Kernel.....	40
3.2.4.1	Depthwise 2D Convolution Kernel	41
3.2.4.2	Pointwise 2D Convolution Kernel	46
3.3	Activation Kernels	50
3.3.1	Sigmoid	50

3.3.2	Tanh	52
3.3.3	Rectifier Linear Unit (ReLU)	54
3.3.4	Softmax	57
3.3.5	Activation Min Max	58
3.3.6	Hard Swish	60
3.3.7	Parametric ReLU (PReLU)	61
3.4	Pooling Kernels	63
3.4.1	Average Pool Kernel	63
3.4.2	Max Pool Kernel	66
3.5	Fully connected Layer	70
3.5.1	Fully Connected Kernel	70
3.6	Basic Operations and Miscellaneous Kernels	73
3.6.1	Interpolation Kernel	73
3.6.2	Dot Product Kernels	74
3.6.3	Elementwise Quantize Kernel	75
3.6.4	Basic Kernels	77
3.6.5	Basic Kernels with Broadcasting	80
3.6.6	Elementwise Comparison Kernels	81
3.6.7	Elementwise Logical Kernels	83
3.6.8	L2 Normalization Kernel	85
3.6.9	Reduce Kernels	86
4.	HiFi 5 NN Library – Layers	90
4.1	GRU Layer	90
4.1.1	GRU Layer Specification	90
4.1.2	Error Codes Specific to GRU	91
4.1.3	API Functions Specific to GRU	92
4.1.3.1	Query Functions	92
4.1.3.2	Initialization Stage	94
4.1.3.3	Execution Stage	95
4.1.4	Structures Specific to GRU	98
4.1.5	Enums Specific to GRU	99
4.2	LSTM Layer	101
4.2.1	LSTM Layer Specification	101
4.2.2	Error Codes Specific to LSTM	101
4.2.3	API Functions Specific to LSTM	103
4.2.3.1	Query Functions	103
4.2.3.2	Initialization Stage	105
4.2.3.3	Execution Stage	106
4.2.4	Structures Specific to LSTM	110
4.2.5	Enums Specific to LSTM	111
4.3	CNN Layer	113
4.3.1	CNN Layer Specification	113
4.3.2	Error Codes Specific to CNN	113
4.3.3	API Functions Specific to CNN	114

4.3.3.1	Query Functions.....	114
4.3.3.2	Initialization Stage.....	117
4.3.3.3	Execution Stage.....	119
4.3.4	Structures Specific to CNN.....	122
4.3.5	Enums Specific to CNN.....	123
5.	Introduction to the Example Testbench.....	125
5.1	Making the Library.....	125
5.2	Making the Executable.....	125
5.3	Sample Testbench for Matrix X Vector Multiplication Kernels.....	126
5.3.1	Usage.....	126
5.4	Sample Testbench for Convolution Kernels.....	127
5.4.1	Usage.....	127
5.5	Sample Testbench for Activation Kernels.....	130
5.5.1	Usage.....	130
5.6	Sample Testbench for Pooling Kernels.....	131
5.6.1	Usage.....	131
5.7	Sample Testbench for Basic Operations Kernels.....	133
5.7.1	Usage.....	133
5.8	Sample Testbench for Normalization Kernels.....	135
5.8.1	Usage.....	135
5.9	Sample Testbench for GRU Layer.....	136
5.9.1	Usage.....	136
5.10	Sample Testbench for LSTM Layer.....	136
5.10.1	Usage.....	137
5.11	Sample Testbench for CNN Layer.....	137
5.11.1	Usage.....	138
5.12	Sample Testbench of TF Micro Lite Example.....	140
5.12.1	Usage.....	140
6.	References.....	141

Figures

Figure 2-1 HiFi NN Layer Interfaces	5
Figure 2-2 Matrix and Cube (SHAPE_CUBE_DWH_T) Shape Representation	6
Figure 2-3 NN Layer Flow Overview	9

Tables

Table 2-1 Library Identification Functions	10
Table 4-1 GRU Get Persistent Size Function	92
Table 4-2 GRU Get Scratch Size Function	93
Table 4-3 GRU Init Function	94
Table 4-4 GRU Execution Function	95
Table 4-5 GRU Set Parameter Function Details	96
Table 4-6 GRU Get Parameter Function Details	97
Table 4-7 GRU Config Structure xa_nnlib_gru_init_config_t	98
Table 4-8 xa_nnlib_gru_weights_t Parameter Type	98
Table 4-9 xa_nnlib_gru_biases_t Parameter Type	99
Table 4-10 Enum xa_nnlib_gru_precision_t	99
Table 4-11 GRU Specific Parameters	100
Table 4-12 LSTM Get Persistent Size Function	103
Table 4-13 LSTM Get Scratch Size Function	104
Table 4-14 LSTM Init Function	105
Table 4-15 LSTM Execution Function	106
Table 4-16 LSTM Set Parameter Function Details	108
Table 4-17 LSTM Get Parameter Function Details	109
Table 4-18 LSTM Config Structure xa_nnlib_lstm_init_config_t	110
Table 4-19 xa_nnlib_lstm_weights_t Parameter Type	110
Table 4-20 xa_nnlib_lstm_biases_t Parameter Type	111
Table 4-21 Enum xa_nnlib_lstm_precision_t	111
Table 4-22 LSTM Specific Parameters	112
Table 4-23 CNN Get Persistent Size Function	114
Table 4-24 CNN Get Scratch Size Function	115
Table 4-25 CNN Init Function	117

Table 4-26 CNN Execution Function	119
Table 4-27 CNN Set Parameter Function Details.....	120
Table 4-28 CNN Get Parameter Function Details	121
Table 4-29 CNN Config Structure xa_nnlib_cnn_init_config_t	122
Table 4-30 Enum xa_nnlib_cnn_precision_t	123
Table 4-31 Enum xa_nnlib_cnn_algo_t	123
Table 4-32 CNN Specific Parameters.....	124

Document Change History

Version	Changes
1.0	Initial version
1.1	Added quantized 8 bit variants for depthwise convolution, fully connected and softmax
1.2	Added quantized 8 bit kernels for SVDF support and for standard convolution, average pooling, and quantization.
1.3	Added quantized 8 bit kernels (asymmetric int8) for pointwise convolution, max pooling, elementwise addition and multiplication.
1.4	Added description of 8 bit kernels (asymmetric int8) for elementwise compare operations, elementwise subtraction, maximum, minimum, tanh, prelu and hardswish.
1.5	<p>Added description of logical operation kernels (boolean 1-byte) for elementwise logical_and, logical_or, logical_not operations.</p> <p>Added per channel quantized sym8s kernel, asym8s input, asym8s output – dilated 2D convolution (without stride support).</p> <p>Added description of L2 Normalization, reduce mean, reduce max kernels (asymmetric int8).</p> <p>Added description of broadcast variants of elementwise maximum, minimum kernels.</p>

1. Introduction to the HiFi 5 NN Library

The HiFi 5 Neural Network (NN) Library is a HiFi-optimized implementation of various NN layers and low level NN kernels. The library is designed with speech and audio neural network domain focus. The low level NN kernels are HiFi-optimized building blocks for NN layer implementation with a generic and simple interface. The NN layers are built using low level kernels and accept input in the form of 'shapes'¹ (up to four dimensions) and produce the output, also in the form of shapes. The layers use the weights or coefficients and biases stored 'externally'² for their operation. The shape of the input, output, weights and biases are as per the layer's design.

This guide refers to the HiFi 5 NN Library as HiFi NN Library, NN layers simply as layers and low level NN kernels as low-level kernels. The current version of the library implements GRU, LSTM (forward path), and CNN layers. It also implements matrix vector multiply, activation, pooling, and convolution functions as low-level kernels.

Note This version of the library supports HiFi 5 DSPs with the NN Extension enabled. The SP-VFPU (Single Precision Vector Floating Point Unit) is optional. The library can be compiled for HiFi 5 DSPs with or without the SP-VFPU enabled.

1.1 Organization of the HiFi 5 NN Library Package

The HiFi NN Library package includes the HiFi NN library containing all layers and low-level kernels implementations, and a set of sample test applications (for layers and low-level kernels).

The HiFi NN library implements a set of NN layers. The application can instantiate these layers and connect inputs and outputs across the layers to form a Neural Network system.

The HiFi NN library also provides a set of low level NN kernels. The application can use these kernels to implement or optimize performance of other NN layers.

The sample test applications implement a file-based application to test an instance of a layer or low level NN kernels for the given specification using pre-generated input, weight or coefficients and bias shapes stored in files in raw binary format.

¹ Refer to Section 2.1 Shape

² Refer to Section 2.2.3 Weights and Biases Memory

1.1.1 Document Overview

This document covers all the information required to integrate the HiFi NN Library into a Neural Network system. All the layers implement “HiFi NN layer APIs”, which is generic and explained in Section 2. The low level NN kernels are explained in Section 3. The APIs for each layer are described in Section 3.6.2. Section 5 provides details about available sample testbenches. References are listed in Section 6.

1.2 HiFi 5 NN Library Specification

The current version of the HiFi NN Library provides the following HiFi-optimized low-level kernels and layer implementations.

1.2.1 Low Level Kernels

- Matrix-vector multiplication kernels
- Convolution kernels
- Activation kernels
- Pooling kernels
- Basic operations kernels

These kernels support fixed point 8 bit, 16 bit, and single precision floating point data types for weights or coefficients, biases, input, and output. Refer to Section 3 for details.

Additionally, 8-bit and 16-bit quantized datatypes as defined in TensorFlow (TF), TensorFlow Lite for Microcontrollers (TFLM) are also supported for select kernels [3]. These datatypes use 8-bit/16-bit quantized values (asym8u – asymmetric 8-bit unsigned, asym8s – asymmetric 8-bit signed, sym8s – symmetric 8-bit signed) for weights or coefficients, input, and output. Biases are 32-bit quantized values.

8-bit quantized types are either unsigned (0, 255) or signed (-128, 127) 8-bit integer with 3 additional parameters.

Three numbers are associated with a quantized 8-bit value that can be used to convert the 8-bit integer to the real value and vice versa. These numbers are:

- Shift: an integer value indicating the amount of shift. If the value is positive, it is left shift and if negative, it is right shift
- Multiplier: a 32 bit (Q31) fixed point value greater than zero.
- Zero point: a 32 bit integer, in range [0, 255] for unsigned type, in range [-128, 127] for signed type.

The formula is:

$$\text{real_value} = (\text{quantized_value} - \text{zero_point}) * 2^{(\text{shift})} * \text{multiplier}$$

The ‘sym8s’ type is symmetrical around 0, this means that quantized values are between -127 to 127 and zero point is 0, so all the calculation required due to zero point is avoided.

To match the asym8u/asym8s/sym8s APIs with Tensorflow, we define zero point as zero_bias in the NN library APIs. The zero_bias is an integer value having range asym8u - [0, 255], asym8s – [-128, 127] (or asym8u - [-255, 0], asym8s – [-127, 128] in case of the reverse operation depending on the corresponding Tensorflow kernel).

In addition to the quantized 8-bit datatypes, a similar 16-bit quantized datatype (asym16s) is used for few kernels. The zero_bias for asym16s datatype is an integer value having range – [-32768, 32767].

1.2.2 Layers

- GRU layer (8x16, 16x16 precision)
- LSTM (forward path) layer (8x16, 16x16 precision)
- CNN layer (8x8, 8x16, 16x16, and float32xfloat32 precision)

Note: MxN precision above denotes (weights or coefficients) x (input, output, bias) precision. Refer to Section 3.6.2 for details.

1.2.3 Support for TensorFlow Lite Micro Operators

The HiFi 5 NN Library low level kernels can be used to implement the following operators of TensorFlow Lite Micro:

No.	Operator	Float32 Datatype Support	Uint8 (asymmetric quantized uint8) Datatype Support	Int8 (quantized int8) Datatype Support	Boolean (1 Byte) Datatype Support
1	FULLY_CONNECTED		Yes	Yes	
2	MAX_POOL_2D	Yes		Yes	
3	SOFTMAX		Yes	Yes	
4	LOGISTIC	Yes		Yes	
5	SVDF			Yes	
6	CONV_2D	Yes	Yes	Yes	
7	DEPTHWISE_CONV_2D	Yes	Yes	Yes	
8	AVERAGE_POOL_2D	Yes	Yes	Yes	
9	FLOOR	Yes			
10	RELU	Yes	Yes	Yes	
11	RELU6	Yes	Yes	Yes	
12	ADD			Yes	
13	MUL			Yes	
14	QUANTIZE			Yes	
15	EQUAL			Yes	
16	GREATER			Yes	
17	GREATEREQUAL			Yes	

18	HARDSWISH			Yes	
19	LESS			Yes	
20	LESSEQUAL			Yes	
21	MAXIMUM			Yes	
22	MINIMUM			Yes	
23	NOTEQUAL			Yes	
24	PRELU			Yes	
25	SUB			Yes	
26	TANH			Yes	
27	LOGICALAND				Yes
28	LOGICALOR				Yes
29	LOGICALNOT				Yes
30	L2 NORM			Yes	
31	MEAN			Yes	
32	REDUCEMAX			Yes	

2. Generic HiFi NN Layer API

Note This section explains an API standard that is evolving. The APIs may undergo some changes in future versions.

This section describes the API that is common to all the HiFi NN layers. The API facilitates any layer instance that works in the overall method shown in Figure 2-1.

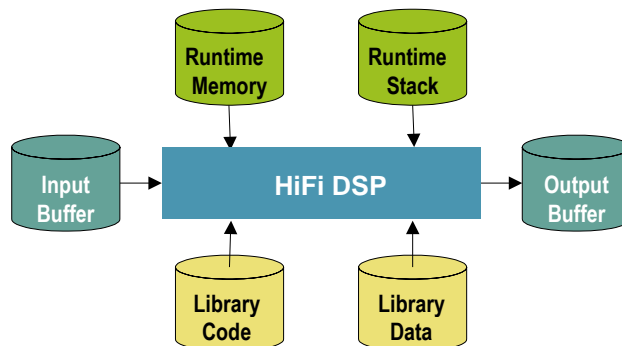


Figure 2-1 HiFi NN Layer Interfaces

All the buffers, input, output, weights and biases are described as shapes. Section 2.1 explains the shape structure.

Section 2.2 discusses all the types of runtime memory required by the layer instances. There is no state information held in static memory, therefore a single thread can perform time division processing of multiple layer instances. Additionally, multiple threads can perform concurrent layer instance processing.

The output from one instance can be fed as input to the next instance if the precision and the dimension matches.

The data types, structures, and error codes explained in this section are declared/defined in `xa_nnlib_standard.h`. By default, the API header file of each layer will include this header file. The application need not include this file.

2.1 Shape

The shapes are used to describe any buffer used in the NN library. The structure `xa_nnlib_shape_t` is defined in `xa_nnlib_standard.h`. The shape can be vector, matrix, or cube.

- Vector is a one-dimensional shape specified by length.
- Matrix is a two-dimensional shape specified by rows, columns, and row_offset. This assumes that the elements in a row are stored at consecutive addresses in memory.

- Cube is a three-dimensional shape specified by height, width, depth, height_offset, width_offset, and depth offset. Cube supports the following shape types:
 - `SHAPE_CUBE_DWH_T`
This assumes that elements are stored in depth, width, and height order; that is, elements with the same height and width indices are stored consecutively.
 - `SHAPE_CUBE_WHD_T`
This assumes that elements are stored in width, height, and depth order; that is, elements with the same height and depth are stored consecutively.

Figure 2-2 shows the dimension variables of matrix and cube shapes.

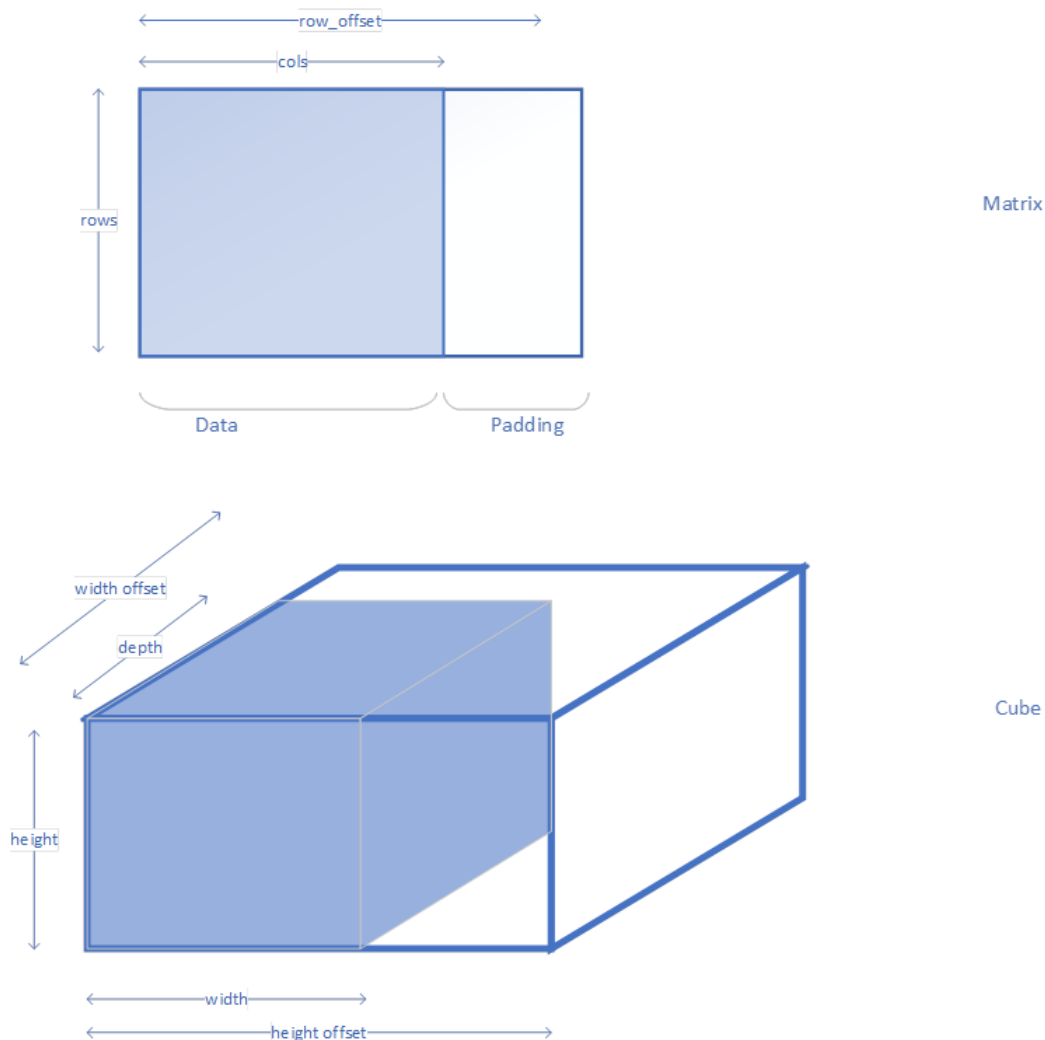


Figure 2-2 Matrix and Cube (`SHAPE_CUBE_DWH_T`) Shape Representation

2.2 Memory Management

The HiFi NN layer API supports a flexible memory scheme and a simple interface that eases the integration into the final application. The API allows the layers to request the required memory for their operations during runtime.

The runtime memory requirement consists primarily of the scratch and persistent memory. The components also require an input buffer and output buffer for the passing of data into and out of the layer.

2.2.1 API Handle / Persistent Memory

The layer API stores persistent state information in a structure that is referenced via an opaque handle. The handle is passed by the application for each API call. This object contains all state and history information that is maintained from one-layer frame invocation to the next within the same thread or instance. The layers expect that the contents of the persistent memory be unchanged by the system apart from the layer itself for the complete lifetime of the layer.

2.2.2 Scratch Memory

This is the temporary buffer used by the layer during a single frame processing call. The contents of this memory region should not be changed if the actual layer execution process is active; that is, if the thread running the layer is inside any API call. This region can be used freely by the system between successive calls to the layer.

2.2.3 Weights and Biases Memory

The weights or coefficients and biases should be managed by the application, and memory should not be requested by the API. If the design requires DMA access from or to the internal memory for better performance, a ping-pong or circular buffer is allocated as part of the scratch into which the weights, biases, input, and output are copied using DMA. If required, these memories can also be persistent.

2.2.4 Input Buffer

This is the buffer from which the layer reads the input. This buffer must be made available for the layer before its execution call. The input buffer should have an associated shape information to describe the input data format. The input buffer pointer can be changed by the application between calls to the layer, but shape information cannot be changed. This allows the layer to read directly from the output of another layer.

2.2.5 Output Buffer

This is the buffer to which the layer writes the output. This buffer must be made available for the layer before its execution call. The output buffer should have an associated shape information to which the layer can describe the output data format. The output buffer pointer can be changed by the application between calls to the layer. This allows the layer to write directly to the input of another layer.

2.3 Generic API Errors

Layer API functions return an error code of type `Int32`, which is of type `signed int`. The format of the error codes is defined in the following table.

31	30 – 27	26–12	11 – 7	6 – 0
Fatal	Class	Reserved	Component	Sub code

The errors that can be returned from the API are subdivided into those that are fatal, which require resetting the layer; and those that are nonfatal and are provided for information to the application.

The class of an error can be API, Config, or Execution. The API category errors are concerned with the incorrect use of the API. The Config errors are produced when the layer parameters are incorrect or outside the supported usage. The Execution errors are returned after a call to the main process and indicate situations that have arisen due to the input data.

2.3.1 Common API Errors

The following errors are fatal and should not be encountered during normal application operation. They signal that a serious error has occurred in the application that is calling the layer.

- `XA_NNLIB_FATAL_MEM_ALLOC`
At least one of the pointers passed into the API function is `NULL`.
- `XA_NNLIB_FATAL_MEM_ALIGN`
At least one of the pointers passed into the API function is not properly aligned.
- `XA_NNLIB_FATAL_INVALID_SHAPE`
At least one of the shapes passed to the API function is invalid.

2.4 C Language API

An overview of the NN layer flow is shown in Figure 2-3. The NN layer API consists of query, initialization, and execution functions.

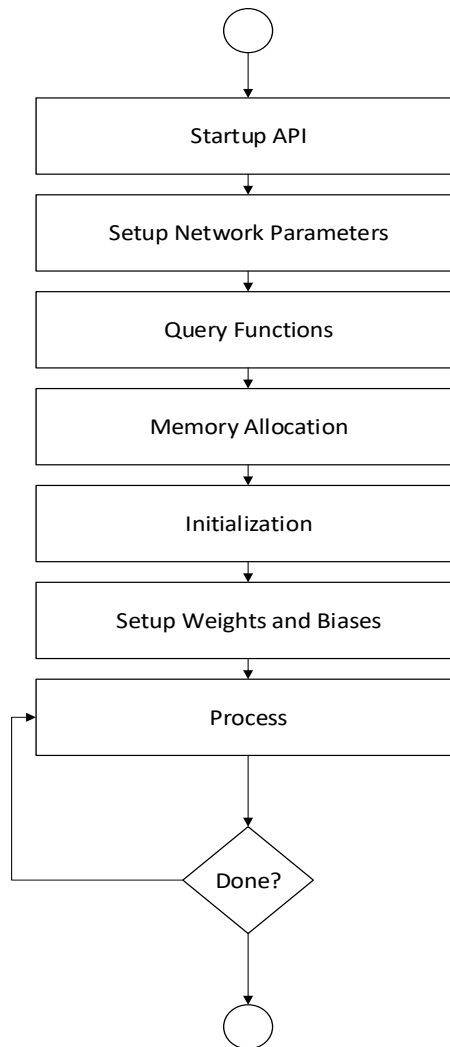


Figure 2-3 NN Layer Flow Overview

2.4.1 Startup Functions

The API startup functions shown in Table 2-1 get the various identification strings from the component library. They are for information only and their usage is optional. These functions do not take any input arguments and return `const char *`.

Table 2-1 Library Identification Functions

Function	Description
<code>xa_nnlib_get_lib_name_string</code>	Get the name of the library.
<code>xa_nnlib_get_lib_version_string</code>	Get the version of the library.
<code>xa_nnlib_get_lib_api_version_string</code>	Get the version of the API.

Example

```
const char *name = xa_nnlib_get_lib_name_string();  
const char *ver = xa_nnlib_get_lib_version_string();  
const char *aver = xa_nnlib_get_lib_api_version_string();
```

Errors

- None

2.4.2 Query Functions

The query functions are used in the startup and the memory allocation stages to obtain information about the memory requirements of the library.

Following is the naming convention for query functions:

```
xa_nnlb_<layer>_get_{persistent | scratch}_<placement>
```

Where:

<layer> indicates the module name (such as gru).

<placement> specifies fast or slow.

2.4.3 Initialization Functions

The initialization functions are used to reset the layer to its initial state. Because the layers are fully reentrant, the application can initialize the layer multiple times.

Following is the naming convention for initialization functions:

```
xa_nnlb_<layer>_init
```

2.4.4 Execution Functions

The execution functions are used to generate the output shape by processing one input shape.

Following is the naming convention for execution functions:

```
xa_nnlb_<layer>_process
```

3. HiFi 5 NN Library – Low-Level Kernels

This section explains the low-level kernels provided in the NN library. All the low-level kernels have a generic, simple interface.

The NN library is a single archive containing all low-level kernels and layers implementations. The following sections explain each low-level kernel in detail.

3.1 *Matrix X Vector Multiplication Kernels*

3.1.1 Matrix X Vector Kernels

Description

These kernels perform the dual matXvec operation with bias addition; that is, $z = \text{mat1} * \text{vec1} + \text{mat2} * \text{vec2} + \text{bias}$. The column dimension of `mat1` must match the row dimension of `vec1` and similarly for `mat2`, `vec2`. Bias and resulting output vector `z` have as many rows as `mat1` and `mat2`.

`bias_shift` and `acc_shift` arguments are provided in the kernel API to adjust Q format of bias and output, respectively. Both `bias_shift` and `acc_shift` can be either positive or negative, where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as matXvec multiplication – accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note, `acc_shift` and `bias_shift` arguments are not relevant in case of floating-point kernels and quantized 8-bit kernels.

`row_stride1` and `row_stride2` arguments are provided in kernel API for row offsets of `mat1` and `mat2`, respectively. Note, input matrices are expected to be appropriately padded in case of `row_stride > cols`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

The arguments, `mat1_zero_bias`, `mat2_zero_bias`, `vec1_zero_bias`, `vec2_zero_bias`, are provided to convert the quantized 8-bit inputs into their real values and perform matXvec operation. The `out_zero_bias`, `out_multiplier` and `out_shift` values are used to quantize real values of output back to 8-bit.

Function variants available are `xa_nn_matXvec_[p]x[q]_[r]`, where:

- `[p]`: Matrix precision in bits
- `[q]`: Vector precision in bits

- [r]: Output precision in bits

Precision

There are twelve variants available:

Type	Description
16x16_16	16-bit matrix inputs, 16-bit vector inputs, 16-bit output
16x16_32	16-bit matrix inputs, 16-bit vector inputs, 32-bit output
16x16_64	16-bit matrix inputs, 16-bit vector inputs, 64-bit output
8x16_16	8-bit matrix inputs, 16-bit vector inputs, 16-bit output
8x16_32	8-bit matrix inputs, 16-bit vector inputs, 32-bit output
8x16_64	8-bit matrix inputs, 16-bit vector inputs, 64-bit output
8x8_8	8-bit matrix inputs, 8-bit vector inputs, 8-bit output
8x8_16	8-bit matrix inputs, 8-bit vector inputs, 16-bit output
8x8_32	8-bit matrix inputs, 8-bit vector inputs, 32-bit output
f32xf32_f32	float32 matrix inputs, float32 vector inputs, float32 output
asym8uxasym8u_asym8u	asym8u matrix inputs, asym8u vector inputs, asym8u output
sym8sxsasym8s_asym8s	sym8s matrix inputs, asym8s vector inputs, asym8s output

Algorithm

$$z_n = 2^{acc-shift} \left(\sum_{m=0}^{cols1-1} mat1_{n,m} \cdot vec1_m + \sum_{m=0}^{cols2-1} mat2_{n,m} \cdot vec2_m + 2^{bias-shift} bias_n \right)$$

For floating-point and quantized 8-bit routines, `acc_shift=0` and `bias_shift=0`.

Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

Prototype

```
WORD32 xa_nn_matXvec_16x16_16
(WORD16 * p_out,          WORD16 * p_mat1,          WORD16 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,          WORD16 * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,
 WORD32 acc_shift,        WORD32 bias_shift);
WORD32 xa_nn_matXvec_16x16_32
(WORD32 * p_out,          WORD16 * p_mat1,          WORD16 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,          WORD16 * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,
 WORD32 acc_shift,        WORD32 bias_shift);
WORD32 xa_nn_matXvec_16x16_64
(WORD64 * p_out,          WORD16 * p_mat1,          WORD16 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,          WORD16 * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,
 WORD32 acc_shift,        WORD32 bias_shift);
WORD32 xa_nn_matXvec_8x16_16
(WORD16 * p_out,          WORD8 * p_mat1,          WORD8 * p_mat2,
```

```

WORD16 * p_vec1,          WORD16 * p_vec2,          WORD16 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,
WORD32 acc_shift,         WORD32 bias_shift);
WORD32 xa_nn_matXvec_8x16_32
(WORD32 * p_out,          WORD8 * p_mat1,          WORD8 * p_mat2,
WORD16 * p_vec1,          WORD16 * p_vec2,          WORD16 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,
WORD32 acc_shift,         WORD32 bias_shift);
WORD32 xa_nn_matXvec_8x16_64
(WORD64 * p_out,          WORD8 * p_mat1,          WORD8 * p_mat2,
WORD16 * p_vec1,          WORD16 * p_vec2,          WORD16 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,
WORD32 acc_shift,         WORD32 bias_shift);
WORD32 xa_nn_matXvec_8x8_8
(WORD8 * p_out,           WORD8 * p_mat1,          WORD8 * p_mat2,
WORD8 * p_vec1,           WORD8 * p_vec2,          WORD8 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,
WORD32 acc_shift,         WORD32 bias_shift);
WORD32 xa_nn_matXvec_8x8_16
(WORD16 * p_out,          WORD8 * p_mat1,          WORD8 * p_mat2,
WORD8 * p_vec1,           WORD8 * p_vec2,          WORD8 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,
WORD32 acc_shift,         WORD32 bias_shift);
WORD32 xa_nn_matXvec_8x8_32
(WORD32 * p_out,          WORD8 * p_mat1,          WORD8 * p_mat2,
WORD8 * p_vec1,           WORD8 * p_vec2,          WORD8 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,
WORD32 acc_shift,         WORD32 bias_shift);
WORD32 xa_nn_matXvec_f32xf32_f32
(FLOAT32 * p_out,         const FLOAT32 * p_mat1, const FLOAT32 * p_mat2,
const FLOAT32 * p_vec1,   const FLOAT32 * p_vec2, const FLOAT32 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2);
WORD32 xa_nn_matXvec_asym8uxasym8u_asym8u
(UWORD8 * p_out,          const UWORD8 * p_mat1, const UWORD8 * p_mat2,
const UWORD8 * p_vec1,   const UWORD8 * p_vec2, const WORD32 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,   WORD32 mat1_zero_bias,
WORD32 mat2_zero_bias,    WORD32 vec1_zero_bias, WORD32 vec2_zero_bias,
WORD32 out_multiplier,    WORD32 out_shift,     WORD32 out_zero_bias);
WORD32 xa_nn_matXvec_sym8sxasym8s_asym8s
(WORD8 * p_out,           const WORD8 * p_mat1, const WORD8 * p_mat2,
const WORD8 * p_vec1,     const WORD8 * p_vec2, const WORD32 * p_bias,
WORD32 rows,              WORD32 cols1,          WORD32 cols2,
WORD32 row_stride1,       WORD32 row_stride2,   WORD32 vec1_zero_bias,
WORD32 vec2_zero_bias,    WORD32 out_multiplier, WORD32 out_shift,
WORD32 out_zero_bias);

```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_mat1	rows*cols1	Input matrix 1, fixed, floating point, asym8u or sym8s
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_mat2	rows*cols2	Input matrix 2, fixed, floating point, asym8u or sym8s
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_vec1	cols1*1	Input vector 1, fixed, floating point, asym8u or asym8s
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_vec2	cols2*1	Input vector 2, fixed, floating point, asym8u or asym8s
WORD16 *, WORD8 *, const WORD32 *, const FLOAT32 *	p_bias	rows*1	Bias vector, fixed or floating point
WORD32	rows		Number of rows in matrix 1, 2 and bias
WORD32	cols1		Number of columns in matrix 1 and rows in vector 1
WORD32	cols2		Number of columns in matrix 2 and rows in vector 2
WORD32	row_stride1		Row offset of matrix 1
WORD32	row_stride2		Row offset of matrix 2
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	mat1_zero_bias		Zero offset of matrix 1
WORD32	mat2_zero_bias		Zero offset of matrix 2
WORD32	vec1_zero_bias		Zero offset of vector 1
WORD32	vec2_zero_bias		Zero offset of vector 2
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
Output			

Type	Name	Size	Description
WORD8 *, UWORD8 *, WORD16 *, WORD32 *, WORD64 *, FLOAT32 *	p_out	rows*1	Output, fixed, floating point, asym8u or asym8s.

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
row_stride1, row_stride2, cols1, cols2	row_stride1 >= cols1 row_stride2 >= cols2
p_mat1, p_mat2, p_vec1, p_vec2, p_bias, p_out	Aligned on <size of one element> boundary Should not overlap
p_mat1, p_vec1, p_out	Cannot be NULL
acc_shift, bias_shift, out_shift	{-31, ..., 31}
mat1_zero_bias, mat2_zero_bias, vec1_zero_bias, vec2_zero_bias	{-255....., 0} for asym8u, {-127....., 128} for asym8s
out_multiplier	Greater than 0
out_zero_bias	{0,.....,255} if out type is asym8u, {-128.....,127} if out type is asym8s

3.1.2 Fused (Activation) Matrix X Vector Kernels

Description

These kernels perform the fused dual matXvec operation with an activation function i.e. $z = \text{activation}(\text{mat1} * \text{vec1} + \text{mat2} * \text{vec2} + \text{bias})$. The column dimension of `mat1` must match the row dimension of `vec1` and similarly for `mat2`, `vec2`. Bias and resulting output vector `z` have as many rows as `mat1` and `mat2`.

Intermediate output of $(\text{mat1} * \text{vec1} + \text{mat2} * \text{vec2} + \text{bias})$ is stored in temporary memory provided by the `p_scratch` argument to kernel API. Activation function is applied on this intermediate output to get final output. Note, for fixed point kernels, the activation function always takes input in Q6.25 format.

`bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and intermediate output respectively. Both `bias_shift` and `acc_shift` can be either positive or negative, where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as `matXvec` multiplication – accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the intermediate output in Q6.25 format.

Note: `acc_shift` and `bias_shift` are not relevant in case of floating point kernels.

`row_stride1` and `row_stride2` arguments are provided in kernel API for row offsets of `mat1` and `mat2` respectively. Note, input matrices are expected to be appropriately padded in case of `row_stride > cols`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

Function variants available are `xa_nn_matXvec_[p]x[q]_[r]_<activation>`, where:

- `[p]`: Matrix precision in bits
- `[q]`: Vector precision in bits
- `[r]`: Output precision in bits
- `<activation>`: activation tag `'tanh'` or `'sigmoid'`

Precision

There are eight variants available:

Type	Description
16x16_16_tanh	16-bit matrix inputs, 16-bit vector inputs, 16-bit output with tanh activation function
16x16_16_sigmoid	16-bit matrix inputs, 16-bit vector inputs, 16-bit output with sigmoid activation function
8x16_16_tanh	8-bit matrix inputs, 16-bit vector inputs, 16-bit output with tanh activation function
8x16_16_sigmoid	8-bit matrix inputs, 16-bit vector inputs, 16-bit output with sigmoid activation function
8x8_8_tanh	8-bit matrix inputs, 8-bit vector inputs, 8-bit output with tanh activation
8x8_8_sigmoid	8-bit matrix inputs, 8-bit vector inputs, 8-bit output with sigmoid activation
f32xf32_f32_tanh	float32 matrix inputs, float32 vector inputs, float32 output with tanh activation
f32xf32_f32_sigmoid	float32 matrix inputs, float32 vector inputs, float32 output with sigmoid activation

Algorithm

$$z_n = \text{activation} \left(2^{\text{acc_shift}} \left(\sum_{m=0}^{\text{cols1}-1} \text{mat1}_{n,m} \cdot \text{vec1}_m + \sum_{m=0}^{\text{cols2}-1} \text{mat2}_{n,m} \cdot \text{vec2}_m + 2^{\text{bias_shift}} \text{bias}_n \right) \right), \quad n = 0, \dots, \text{rows} - 1$$

In case of floating point routine, `acc_shift=0` and `bias_shift=0`.

Thus, $2^{acc_shift} = 2^{bias_shift} = 1$

activation is tanh or sigmoid

Prototype

```
WORD32 xa_nn_matXvec_16x16_16_tanh
(WORD16 * p_out,          WORD16 * p_mat1,          WORD16 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,         VOID * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 bias_precision,    VOID * p_scratch);
WORD32 xa_nn_matXvec_16x16_16_sigmoid
(WORD16 * p_out,          WORD16 * p_mat1,          WORD16 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,         VOID * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 bias_precision,    VOID * p_scratch);
WORD32 xa_nn_matXvec_8x16_16_tanh
(WORD16 * p_out,          WORD8 * p_mat1,           WORD8 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,         VOID * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 bias_precision,    VOID * p_scratch);
WORD32 xa_nn_matXvec_8x16_16_sigmoid
(WORD16 * p_out,          WORD8 * p_mat1,           WORD8 * p_mat2,
 WORD16 * p_vec1,         WORD16 * p_vec2,         VOID * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 bias_precision,    VOID * p_scratch);
WORD32 xa_nn_matXvec_8x8_8_tanh
(WORD8 * p_out,           WORD8 * p_mat1,           WORD8 * p_mat2,
 WORD8 * p_vec1,          WORD8 * p_vec2,          VOID * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 bias_precision,    VOID * p_scratch);
WORD32 xa_nn_matXvec_8x8_8_sigmoid
(WORD8 * p_out,           WORD8 * p_mat1,           WORD8 * p_mat2,
 WORD8 * p_vec1,          WORD8 * p_vec2,          VOID * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 bias_precision,    VOID * p_scratch);
WORD32 xa_nn_matXvec_f32xf32_f32_tanh
(FLOAT32 * p_out,         FLOAT32 * p_mat1,        FLOAT32 * p_mat2,
 FLOAT32 * p_vec1,        FLOAT32 * p_vec2,        FLOAT32 * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      FLOAT32 * p_scratch);
WORD32 xa_nn_matXvec_f32xf32_f32_sigmoid
(FLOAT32 * p_out,         FLOAT32 * p_mat1,        FLOAT32 * p_mat2,
 FLOAT32 * p_vec1,        FLOAT32 * p_vec2,        FLOAT32 * p_bias,
 WORD32 rows,             WORD32 cols1,             WORD32 cols2,
 WORD32 row_stride1,      WORD32 row_stride2,      FLOAT32 * p_scratch);
```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, FLOAT32 *	p_mat1	rows*cols1	Input matrix 1, fixed or floating point
WORD16 *, WORD8 *, FLOAT32 *	p_mat2	rows*cols2	Input matrix 2, fixed or floating point
WORD16 *, WORD8 *, FLOAT32 *	p_vec1	cols1*1	Input vector 1, fixed or floating point
WORD16 *, WORD8 *, FLOAT32 *	p_vec2	cols2*1	Input vector 2, fixed or floating point
VOID *, FLOAT32 *	p_bias	rows*1	Bias vector, fixed or floating point
WORD32	rows		Number of rows in matrix 1,2, bias and output
WORD32	cols1		Number of columns in matrix 1 and rows in vector 1
WORD32	cols2		Number of columns in matrix 2 and rows in vector 2
WORD32	row_stride1		Row offset of matrix 1
WORD32	row_stride2		Row offset of matrix 2
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	bias_precision		Precision of bias in bytes
Output			
WORD8 *, WORD16 *, FLOAT32 *	p_out	rows*1	Output, fixed (Q7, Q15) or floating point
Temporary			
VOID *, FLOAT32 *	p_scratch	rows*4	Scratch (temporary) memory pointer

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
row_stride1, row_stride2, cols1, cols2	Multiples of 4 (2 in case of floating point)
p_mat1, p_mat2, p_vec1, p_vec2, p_bias, p_out	Aligned on 16-byte boundary Should not overlap
p_mat1, p_vec1, p_bias, p_out	Cannot be NULL
acc_shift, bias_shift	{-31, ..., 31}
bias_precision	{-1, 8, 16, 32, 64} (-1 in case of floating point)

3.1.3 Matrix X Vector Batch Kernels

Description

These kernels perform the operation of multiplication of a single matrix with a series of vectors along with bias addition; that is, $z_i = \text{mat1} * \text{vec1}_i + \text{bias}$. These kernels can also be viewed as matrix X matrix-transpose multiplication kernels. The column dimension of `mat1` must match the row dimension of vectors in `vec1`. Bias and resulting output vector sequence `z` have as many numbers of rows as `mat1`. `vec1` is a sequence of `vec_count` number of input vectors and bias is added to each resulting vector after multiplication with `mat1`. Thus, output `z` has dimensions `rows*vec_count`. `vec_count` number of input vectors and output vectors are provided as array of pointers arguments to kernel API.

`bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output respectively. Both `bias_shift` and `acc_shift` can be either positive or negative where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as `matXvec` multiplication – accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note: `acc_shift` and `bias_shift` are not relevant in case of floating point kernels.

The `row_stridel` argument is provided in kernel API for row offset of `mat1`. Note, input matrix is expected to be appropriately padded in case of `row_stridel > cols1`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

Function variants available are `xa_nn_matXvec_batch_[p]x[q]_[r]`, where:

- `[p]`: Matrix precision in bits
- `[q]`: Vector precision in bits
- `[r]`: Output precision in bits

Precision

There are four variants available:

Type	Description
16x16_64	16-bit matrix inputs, 16-bit vector inputs, 64-bit output vectors
8x16_64	8-bit matrix inputs, 16-bit vector inputs, 64-bit output vectors
8x8_32	8-bit matrix inputs, 8-bit vector inputs, 32-bit output vectors
f32xf32_f32	float32 matrix inputs, float32 vector inputs, float32 output

Algorithm

$$z_{n,i} = 2^{acc-shift} \left(\sum_{m=0}^{cols1-1} mat1_{n,m} \cdot vec1_{m,i} + 2^{bias-shift} bias_n \right),$$

$$n = 0, \dots, \overline{rows} - 1 ; \quad i = 0, \dots, \overline{vec-count} - 1$$

In case of floating point routine, acc_shift=0 and bias_shift=0.

Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

Prototype

```
WORD32 xa_nn_matXvec_batch_16x16_64
(WORD64 ** p_out,          WORD16 * p_mat1,          WORD16 ** p_vec1,
 WORD16 * p_bias,          WORD32 rows,              WORD32 cols1,
 WORD32 row_stridel,       WORD32 acc_shift,          WORD32 bias_shift,
 WORD32 vec_count);

WORD32 xa_nn_matXvec_batch_8x16_64
(WORD64 ** p_out,          WORD8 * p_mat1,           WORD16 ** p_vec1,
 WORD16 * p_bias,          WORD32 rows,              WORD32 cols1,
 WORD32 row_stridel,       WORD32 acc_shift,          WORD32 bias_shift,
 WORD32 vec_count);

WORD32 xa_nn_matXvec_batch_8x8_32
(WORD32 ** p_out,          WORD8 * p_mat1,           WORD8 ** p_vec1,
 WORD8 * p_bias,          WORD32 rows,              WORD32 cols1,
 WORD32 row_stridel,       WORD32 acc_shift,          WORD32 bias_shift,
 WORD32 vec_count);

WORD32 xa_nn_matXvec_batch_f32xf32_f32
(FLOAT32 ** p_out,         FLOAT32 * p_mat1,         FLOAT32 ** p_vec1,
 FLOAT32 * p_bias,         WORD32 rows,              WORD32 cols1,
 WORD32 row_stridel,       WORD32 vec_count);
```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, FLOAT32 *	p_mat1	rows*cols1	Input matrix, fixed or floating point
WORD16 **, WORD8 **, FLOAT32 **	p_vec1	cols1*vec_count	Input vector pointers, fixed or floating point
WORD16 *, WORD8 *, FLOAT32 *	p_bias	rows*1	Bias vector, fixed or floating point
WORD32	rows		Number of rows in input matrix, bias and output
WORD32	cols1		Number of columns in input matrix and rows in input vector
WORD32	row_stridel		Row offset of input matrix
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	vec_count		Number of input vectors

Type	Name	Size	Description
Output			
WORD32 **, WORD64 **, FLOAT32 **	p_out	rows*vec_count	Output vector pointers, fixed or floating point

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
row_stride1, cols1	Multiples of 4 (2 in case of floating point)
p_mat1, p_vec1, p_bias, p_out	Aligned on 16-byte boundary Should not overlap Cannot be NULL
acc_shift, bias_shift	{-31, ..., 31}

3.1.4 Matrix Multiplication Kernels

Description

These kernels perform the operation of multiplication of a matrix `mat1` with another matrix `mat2` along with bias addition; that is, $z = \text{mat1} * \text{mat2} + \text{bias}$. The first matrix should be stored in row major order and the second matrix should be stored in column major order. The first matrix is of dimensions `rows` x `cols`. The second matrix `mat2` is of dimensions `cols` x `vec_count`. These kernels can also be viewed as a modification of the Matrix X Vector Batch kernels. The column dimension of `mat1` matches the row dimension of `mat2` i.e. the length of each vector in `p_mat2`. Bias and resulting output vector sequence `z` have as many numbers of rows as `mat1`. `mat2` is a sequence of `vec_count` number of input vectors and bias is added to each resulting vector after multiplication with `mat1`. Thus, output `z` has dimensions `rows` * `vec_count`. The arguments `vec_offset` and `out_offset` are offsets to the next vector and output addresses. The argument `out_stride` defines the row offset for the output matrix. For standard matrix multiplication, `vec_offset` should be equal to `cols`, `out_offset` equal to 1 and `out_stride` should be equal to `vec_count` i.e. columns of `mat2`.

The `bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output respectively. Both `bias_shift` and `acc_shift` can be either positive or negative where positive value denotes a left shift and negative value denotes a right shift.

The `bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as multiplication – accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note, the `acc_shift` and `bias_shift` arguments are not relevant in case of quantized 8-bit kernels.

The `row_stride` argument indicates the offset to next row of `mat1`.

The `vec_offset` argument refers to the column offset of `mat2`.

Similarly, the `out_offset` and `out_stride` arguments refer to the column offset and row offset of the output matrix `rows * vec_count` respectively.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

The arguments, `mat1_zero_bias`, `mat2_zero_bias`, are provided to convert the quantized 8-bit inputs into their real values and perform `matXvec` batch operation. The `out_zero_bias`, `out_multiplier` and `out_shift` values are used to quantize real values of output back to quantized 8-bit values.

For the quantized int8 variant, we have per-row quantized input `mat1`.

Function variants available are `xa_nn_matmul_[p]x[q]_[r]`, where:

- `[p]`: Matrix 1 precision in bits
- `[q]`: Matrix 2 precision in bits
- `[r]`: Output precision in bits

Precision

There are three variants available:

Type	Description
<code>8x8_8</code>	8-bit matrix inputs, 8-bit vector inputs, 8-bit output vectors
<code>asym8uxasym8u_asym8u</code>	asym8u matrix inputs, asym8u vector inputs, asym8u output vectors
<code>per_chan_sym8sxasym8s_asym8s</code>	per channel quantized sym8s matrix inputs, asym8s vector inputs, asym8s output vectors

Algorithm

$$z_{n,i} = 2^{acc-shift} \left(\sum_{m=0}^{cols1-1} mat1_{n,m} \cdot mat2_{m,i} + 2^{bias-shift} bias_n \right),$$

$$n = 0, \dots, \overline{rows} - 1 ; \quad i = 0, \dots, \overline{vec-count} - 1$$

In case of quantized 8-bit routines, `acc_shift=0` and `bias_shift=0`.

Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

Prototype

```
WORD32 xa_nn_matmul_8x8_8
(WORD8 * p_out,          WORD8 * p_mat1,          WORD8 * p_mat2,
 WORD8 * p_bias,         WORD32 rows,              WORD32 cols,
```

```

WORD32 row_stride,          WORD32 acc_shift,          WORD32 bias_shift,
WORD32 vec_count,          WORD32 vec_offset,          WORD32 out_offset,
WORD32 out_stride);
WORD32 xa_nn_matmul_asym8uxasym8u_asym8u
(WORD8 * p_out,             const UWORD8 * p_mat1,      const UWORD8 * p_mat2,
 const WORD32 * p_bias,     WORD32 rows,               WORD32 cols,
 WORD32 row_stride,        WORD32 vec_count,          WORD32 vec_offset,
 WORD32 out_offset,        WORD32 out_stride,          WORD32 mat1_zero_bias,
 WORD32 vec1_zero_bias,    WORD32 out_multiplier,      WORD32 out_shift,
 WORD32 out_zero_bias);
WORD32 xa_nn_matmul_per_chan_sym8sxasym8s_asym8s
(WORD8 * p_out,             const WORD8 * p_mat1,       const WORD8 * p_mat2,
 const WORD32 * p_bias,     WORD32 rows,               WORD32 cols,
 WORD32 row_stride,        WORD32 vec_count,          WORD32 vec_offset,
 WORD32 out_offset,        WORD32 out_stride,          WORD32 vec1_zero_bias,
 const WORD32 *p_out_multiplier, const WORD32 p_out_shift,
 WORD32 out_zero_bias);

```

Arguments

Type	Name	Size	Description
Input			
WORD8 *, UWORD8 *, const WORD8 *	p_mat1	rows*cols	Input matrix 1, fixed point, asym8u or sym8s
WORD8 *, UWORD8 *, const WORD8 *	p_mat2	cols * vec_count	Input matrix 2, fixed , asym8u or sym8s
WORD8 *, const WORD32 *	p_bias	rows*1	Bias vector, fixed point
WORD32	rows		Number of rows in input matrix, bias and output
WORD32	cols		Number of columns in input matrix and rows in input vector
WORD32	row_stride		Row offset of input matrix
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	vec_count		Number of vectors (columns) in matrix 2
WORD32	vec_offset		Offset to the next vector address
WORD32	out_offset		Offset to the next output address
WORD32	out_stride		Row offset of output matrix
WORD32	mat1_zero_bias		Zero offset of matrix 1
WORD32	vec1_zero_bias		Zero offset of matrix 2
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
Output			
WORD8 *, UWORD8 *	p_out	rows*vec_ count	Output matrix, fixed-point, floating point or asym8u

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_mat1, p_mat2, p_out	Aligned on (size of one element)-byte boundary Cannot be NULL Should not overlap
p_bias	Aligned on (size of one element)-byte boundary
acc_shift, bias_shift, out_shift	{-31, ..., 31}
vec_count	Greater than 0
vec_offset, out_offset, out_stride	Should not be 0
mat1_zero_bias,	{-255, ..., 0}
vec1_zero_bias	{-255, ..., 0} for asym8u, {-127, ..., 127} for asym8s
out_multiplier	Greater than 0
out_zero_bias	{0, ..., 255} if out type is asym8u, {-128, ..., 127} if out type is asym8s

3.1.5 Matrix X Vector Kernels with Output Stride

Description

These kernels perform a single matXvec operation with bias addition; that is, $z = \text{mat1} * \text{vec1} + \text{bias}$. The column dimension of `mat1` must match the row dimension of `vec1`. Bias and resulting output vector `z` have as many rows as `mat1`.

`row_stride1` is provided in kernel API for row offsets of `mat1`. Note, input matrices are expected to be appropriately padded in case of `row_stride > cols`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

The argument `out_stride` is helpful in storing the output at a given offset.

The argument `vec1_zero_bias` is provided to convert the quantized 8-bit inputs into their real values and perform matXvec operation. The `out_multiplier` and `out_shift` values are used to convert real values of output to 16-bit.

Function variants available are `xa_nn_matXvec_[p]x[q]_[r]`, where:

- `[p]`: Matrix precision in bits

- [q]: Vector precision in bits
- [r]: Output precision in bits

Precision

There is one variant available:

Type	Description
sym8sxasym8s_16	sym8s matrix inputs, asym8s vector inputs, asym8s output

Algorithm

$$z_n = \left(\sum_{m=0}^{cols1-1} mat1_{n,m} \cdot vec1_m + bias_n \right)$$

Prototype

```
WORD32 xa_nn_matXvec_out_stride_sym8sxasym8s_16
(WORD16 * p_out,      const WORD8 * p_mat1,  const WORD8 * p_vec1,
 const WORD32 * p_bias, WORD32 rows,        WORD32 cols1,
 WORD32 row_stride1,  WORD32 out_stride,     WORD32 vec1_zero_bias,
 WORD32 out_multiplier, WORD32 out_shift);
```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *	p_mat1	rows*cols1	Input matrix, sym8s
const WORD8 *	p_vec1	cols1*1	Input vector, asym8s
const WORD32 *	p_bias	rows*1	Bias vector
WORD32	rows		Number of rows in matrix and number of elements in bias
WORD32	cols1		Number of columns in matrix and elements in vector
WORD32	row_stride1		Row offset of matrix
WORD32	out_stride		Row offset of output
WORD32	vec1_zero_bias		Zero offset of vector
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
Output			
WORD16 *	p_out	rows*1	Output, 16-bit

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
row_stridel, cols1	row_stride1 >= cols1
p_mat1, p_vec1, p_bias, p_out	Aligned on <size of one element> boundary Should not overlap
p_mat1, p_vec1, p_out	Cannot be NULL
out_shift	{-31, ..., 31}
vec1_zero_bias	{-127....., 128} for asym8s
out_multiplier	Greater than 0

3.2 Convolution Kernels

3.2.1 Standard 2D Convolution Kernel

Description

These kernels perform the 2D convolution operation as $z = \text{inp}(\ast)\text{kernel} + \text{bias}$. A 3D input cube ($\text{input_height} \times \text{input_width} \times \text{input_channels}$), is convolved with a 3D kernel cube ($\text{kernel_height} \times \text{kernel_width} \times \text{input_channels}$) to produce a 2D convolution output plane ($\text{out_height} \times \text{out_width}$). With out_channels number of such 3D kernels, output cube ($\text{out_height} \times \text{out_width} \times \text{out_channels}$) is produced. The bias having dimension (out_channels) is added after the convolution (one bias value is added to each output channel) to produce the final output.

Note: The depth or channels dimension (input_channels) of input and kernel must be identical for 2D convolution.

bias_shift and acc_shift arguments are provided in kernel API to adjust Q format of bias and output, respectively. Both bias_shift and acc_shift can be either positive or negative where positive value denotes a left shift and negative value denotes a right shift.

bias_shift is the shift in number of bits applied to the bias to make it in the same Q format as convolution - accumulation result. acc_shift is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note: acc_shift and bias_shift are not relevant in case of floating point kernels and quantized 8-bit kernels.

The x_stride and y_stride arguments in kernel API define the step size of the kernel when traversing the input in width and height dimensions respectively.

The x_padding argument defines padding to the left of the input in the width dimension and the y_padding argument defines padding to the top of the input in the height dimension.

The right padding is calculated based on `out_width` as `right_paddding = kernel_width + (out_width - 1) * x_stride - (x_padding + input_width)`.

The bottom padding is calculated based on `out_height` as `bottom_paddding = kernel_height + (out_height - 1) * y_stride - (y_padding + input_height)`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

The kernel is expected to be padded in the depth or channels dimension if the number of `input_channels` is not a multiple of 4 in case of fixed point variants other than the 8x8, `asym8uxasym8u` and `per_chan_sym8sxasym8s` variant, and 2 in case of floating point variant. No padding is needed for 8x8 and quantized 8-bit variants.

These kernels require temporary buffer for convolution computation. This temporary buffer is provided by `p_scratch` argument of kernel API. The size of temporary buffer should be queried using `xa_nn_conv2d_std_getsize()` helper API.

These kernels expect input and kernel cubes in `SHAPE_CUBE_DWH_T` shape type and can produce output cube in either `SHAPE_CUBE_DWH_T` or `SHAPE_CUBE_WHD_T` shape type. The `out_data_format` argument to kernel API controls the output cube shape type.

Function variants available are `xa_nn_conv2d_std_[p]x[q]`, where:

- `[p]`: Kernel precision in bits
- `[q]`: Input precision in bits

Precision

There are six variants available.

Type	Description
16x16	16-bit kernel, 16-bit input, 16-bit output
8x16	8-bit kernel, 16-bit input, 16-bit output
8x8	8-bit kernel, 8-bit input, 8-bit output
f32	float32 kernel, float32 input, float32 output
asym8uxasym8u	asym8u kernel, asym8u input, asym8u output
per_chan_sym8sxasym8s	per channel quantized sym8s kernel, asym8s input, asym8s output

Algorithm

$$z_{h,w,d} = 2^{acc-shift} \left(\sum_{i=0}^{K_H-1} \sum_{j=0}^{K_W-1} \sum_{k=0}^{I_C-1} in_{pad_{(h*y-stride+i),(w*x-stride+j),k}} \cdot ker_{pad_{d,i,j,k}} + 2^{bias-shift} b_d \right)$$

$$h = 0, \dots, \overline{out-height - 1}, w = 0, \dots, \overline{out-width - 1},$$

$d = 0, \dots, \overline{\text{out-channels}} - 1$

In case of floating-point kernels and quantized 8-bit kernels, $\text{acc_shift}=0$ and $\text{bias_shift}=0$.

Thus, $2^{\text{acc_shift}} = 2^{\text{bias_shift}} = 1$

$\text{in}_{\text{pad}}, \text{ker}_{\text{pad}}$ denote the padded p_inp and padded p_ker shapes, respectively.

K_H, K_W, I_C denote kernel_height, kernel_width, and input_channels, respectively.

b denotes the bias shape.

Prototype

```
WORD32 xa_nn_conv2d_std_getsize
(WORD32 input_height,      WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width,     WORD32 y_stride,      WORD32 y_padding,
 WORD32 out_height,       WORD32 input_precision);

WORD32 xa_nn_conv2d_std_16x16
(WORD16 * p_out,          WORD16 * p_inp,      WORD16 * p_ker,
 WORD16 * p_bias,         WORD32 input_height, WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height, WORD32 kernel_width,
 WORD32 out_channels,     WORD32 x_stride,     WORD32 y_stride,
 WORD32 x_padding,        WORD32 y_padding,    WORD32 out_height,
 WORD32 out_width,        WORD32 bias_shift,   WORD32 acc_shift,
 WORD32 out_data_format,  VOID * p_scratch);

WORD32 xa_nn_conv2d_std_8x16
(WORD16 * p_out,          WORD16 * p_inp,      WORD8 * p_ker,
 WORD16 * p_bias,         WORD32 input_height, WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height, WORD32 kernel_width,
 WORD32 out_channels,     WORD32 x_stride,     WORD32 y_stride,
 WORD32 x_padding,        WORD32 y_padding,    WORD32 out_height,
 WORD32 out_width,        WORD32 bias_shift,   WORD32 acc_shift,
 WORD32 out_data_format,  VOID * p_scratch);

WORD32 xa_nn_conv2d_std_8x8
(WORD8 * p_out,           WORD8 * p_inp,       WORD8 * p_ker,
 WORD8 * p_bias,          WORD32 input_height, WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height, WORD32 kernel_width,
 WORD32 out_channels,     WORD32 x_stride,     WORD32 y_stride,
 WORD32 x_padding,        WORD32 y_padding,    WORD32 out_height,
 WORD32 out_width,        WORD32 bias_shift,   WORD32 acc_shift,
 WORD32 out_data_format,  VOID * p_scratch);

WORD32 xa_nn_conv2d_std_f32
(FLOAT32 * p_out,         const FLOAT32 * p_inp, const FLOAT32 * p_ker,
 const FLOAT32 * p_bias,  WORD32 input_height, WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height, WORD32 kernel_width,
 WORD32 out_channels,     WORD32 x_stride,     WORD32 y_stride,
 WORD32 x_padding,        WORD32 y_padding,    WORD32 out_height,
 WORD32 out_width,        WORD32 out_data_format, VOID * p_scratch);

WORD32 xa_nn_conv2d_std_asym8uxasym8u
(UWORD8 * p_out,          const UWORD8 * p_inp,  const UWORD8 * p_ker,
 const WORD32 * p_bias,   WORD32 input_height, WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height, WORD32 kernel_width,
 WORD32 out_channels,     WORD32 x_stride,     WORD32 y_stride,
 WORD32 x_padding,        WORD32 y_padding,    WORD32 out_height,
 WORD32 out_width,        WORD32 input_zero_bias, WORD32 kernel_zero_bias,
```

```

WORD32 out_multiplier,    WORD32 out_shift,    WORD32 out_zero_bias,
WORD32 out_data_format,  VOID * p_scratch);
WORD32 xa_nn_conv2d_std_per_chan_sym8sxasym8s
(WORD8 * p_out,           const WORD8 * p_inp,    const WORD8 * p_ker,
 const WORD32 * p_bias,   WORD32 input_height,   WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height,  WORD32 kernel_width,
 WORD32 out_channels,     WORD32 x_stride,      WORD32 y_stride,
 WORD32 x_padding,        WORD32 y_padding,      WORD32 out_height,
 WORD32 out_width,        WORD32 input_zero_bias, WORD32 * p_out_multiplier,
 WORD32 * p_out_shift,    WORD32 out_zero_bias, WORD32 out_data_format,
 VOID * p_scratch);

```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_inp	input_height* input width* input_channels	Input cube, fixed, floating point, asym8u or asym8s, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_ker	out_channels* (kernel_height * kernel width* input_channels)	Kernel cube, fixed, floating point, asym8u or sym8s, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, FLOAT32 *, const WORD32 *	p_bias	out_channels	Bias vector, fixed or floating point
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Number of input channels
WORD32	kernel_height		Kernel height
WORD32	kernel_width		Kernel width
WORD32	out_channels		Number of output channels
WORD32	x_stride		Horizontal stride over input
WORD32	y_stride		Vertical stride over input
WORD32	x_padding		Left padding width on input
WORD32	y_padding		Top padding height on input
WORD32	out_height		Output height
WORD32	out_width		Output width
WORD32	bias_shift		Shift applied to bias
WORD32	acc_shift		Shift applied to accumulator

WORD32	input_zero_bias		Zero offset of input
WORD32	kernel_zero_bias		Zero offset of kernel
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
WORD32	out_data_format		Output data format 0:SHAPE_CUBE_DWH_T 1:SHAPE_CUBE_WHD_T
VOID *	p_scratch	xa_nn_conv2d_std_getsize()	Scratch memory pointer
Output			
WORD16 *, WORD8 *, FLOAT32 *, UWORD8 *	p_out	(out_height* out_width)* out_channels	Output cube, fixed, floating point, asym8u or asym8s, as per the out_data_format argument.

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_out, p_inp, p_ker, p_bias, p_scratch	Cannot be NULL
	Should not overlap
	Aligned on 16-byte boundary except for quantized 8-bit kernels where only p_scratch is required to be 16-byte aligned.
	For p_scratch - memory size >= size returned by xa_nn_conv2d_std_getsize()
input_height, input_width, input_channels	Greater than or equal to 1
kernel_height	{1, 2, ..., input_height}
kernel_width	{1, 2, ..., input_width}
out_channels	Greater than or equal to 1
x_stride	{1, 2, ..., kernel_width}
y_stride	Greater than or equal to 1
x_padding, y_padding	Greater than or equal to 0
out_height, out_width	Greater than or equal to 1
acc_shift, bias_shift, out_shift	{-31 31} for fixed point and quantized 8-bit APIs
input_zero_bias	{-255,, 0} for asym8u input, {-127,, 128} for asym8s input
kernel_zero_bias	{-255,, 0} for asym8u kernel
out_zero_bias	{0,, 255} for asym8u output, {-128,, 127} for asym8s output
out_multiplier	Greater than 0

out_data_format	Can be 0: SHAPE_CUBE_DWH_T or 1: SHAPE_CUBE_WHD_T
-----------------	--

3.2.2 Standard 2D Convolution Kernel with Dilation

Description

These kernels perform the dilated 2D convolution operation as $z = \text{inp}(\ast)\text{kernel} + \text{bias}$. A 3D input cube (`input_height` x `input_width` x `input_channels`), is convolved with a 3D dilated kernel cube to produce a 2D convolution output plane (`out_height` x `out_width`). With `out_channels` number of such 3D kernels, output cube (`out_height` x `out_width` x `out_channels`) is produced. Before convolution, the 3D kernel cube (`kernel_height` x `kernel_width` x `input_channels`) is dilated by skipping `dilation_height-1` elements in height dimension and `dilation_width-1` elements in width dimension with, `dilation_height`>=1 and/or `dilation_width`>=1. Post dilation, the kernel cube is of size `kernel_height_dilation = kernel_height + (kernel_height-1)*(dilation_height-1)` in height dimension and `kernel_width_dilation = kernel_width + (kernel_width-1)*(dilation_width-1)` in width dimension. The bias having dimension (`out_channels`) is added after the convolution (one bias value is added to each output channel) to produce the final output.

Note: The depth or channels dimension (`input_channels`) of input and kernel must be identical for 2D convolution.

`bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output, respectively. Both `bias_shift` and `acc_shift` can be either positive or negative where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as convolution - accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

The `x_stride` and `y_stride` arguments in kernel API define the step size of the kernel when traversing the input in width and height dimensions respectively. Presently, stride feature is not supported and so, stride values equal to 1 only to be used.

The `x_padding` argument defines padding to the left of the input in the width dimension and the `y_padding` argument defines padding to the top of the input in the height dimension.

The right padding is calculated based on `out_width` as `right_paddding = kernel_width_dilation + (out_width - 1) * x_stride - (x_padding + input_width)`.

The bottom padding is calculated based on `out_height` as `bottom_paddding = kernel_height_dilation + (out_height - 1) * y_stride - (y_padding + input_height)`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

These kernels require temporary buffer for convolution computation. This temporary buffer is provided by `p_scratch` argument of kernel API. The size of temporary buffer should be queried using `xa_nn_dilated_conv2d_std_getsize()` helper API.

These kernels expect input and kernel cubes in `SHAPE_CUBE_DWH_T` shape type and can produce output cube in either `SHAPE_CUBE_DWH_T` or `SHAPE_CUBE_WHD_T` shape type. The `out_data_format` argument to kernel API controls the output cube shape type.

Precision

Type	Description
<code>per_chan_sym8sxasym8s</code>	per channel quantized sym8s kernel, asym8s input, asym8s output

Algorithm

$$\begin{aligned}
 &Z_{h,w,d} \\
 &= 2^{acc-shift} \left(\sum_{i=0}^{K_H-1} \sum_{j=0}^{K_W-1} \sum_{k=0}^{I_C-1} in_{pad}(h*y_stride+i*dilation_height),(w*x_stride+j*dilation_width),k \right. \\
 &\quad \left. \cdot ker_{d,i,j,k} + 2^{bias-shift} b_d \right) \\
 &h = 0, \dots, \overline{out_height - 1}, w = 0, \dots, \overline{out_width - 1}, \\
 &d = 0, \dots, \overline{out_channels - 1}
 \end{aligned}$$

in_{pad} , ker denote the padded `p_inp` and kernel `p_ker` shapes, respectively.

K_H, K_W, I_C denote `kernel_height`, `kernel_width`, and `input_channels`, respectively.

b denotes the `bias` shape.

Prototype

```
WORD32 xa_nn_dilated_conv2d_std_getsize
(WORD32 input_height,      WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width,     WORD32 y_stride,      WORD32 y_padding,
 WORD32 out_height,       WORD32 input_precision, WORD32 dilation_height);
```

```
WORD32 xa_nn_dilated_conv2d_std_per_chan_sym8sxasym8s
(WORD8 * p_out,           const WORD8 * p_inp,   const WORD8 * p_ker,
 const WORD32 * p_bias,   WORD32 input_height,  WORD32 input_width,
 WORD32 input_channels,  WORD32 kernel_height, WORD32 kernel_width,
 WORD32 out_channels,    WORD32 x_stride,      WORD32 y_stride,
 WORD32 x_padding,       WORD32 y_padding,     WORD32 out_height,
 WORD32 out_width,       WORD32 input_zero_bias, WORD32 * p_out_multiplier,
 WORD32 * p_out_shift,    WORD32 out_zero_bias,  WORD32 out_data_format,
 VOID * p_scratch, WORD32 dilation_height, WORD32 dilation_width);
```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_inp	input_height* input_width* input_channels	Input cube, fixed, floating point, asym8u or asym8s, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_ker	out_channels* (kernel_height* kernel_width* input_channels)	Kernel cube, fixed, floating point, asym8u or sym8s, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, FLOAT32 *, const WORD32 *	p_bias	out_channels	Bias vector, fixed or floating point
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Number of input channels
WORD32	kernel_height		Kernel height
WORD32	kernel_width		Kernel width
WORD32	out_channels		Number of output channels
WORD32	x_stride		Horizontal stride over input
WORD32	y_stride		Vertical stride over input
WORD32	x_padding		Left padding width on input
WORD32	y_padding		Top padding height on input
WORD32	out_height		Output height
WORD32	out_width		Output width
WORD32	bias_shift		Shift applied to bias
WORD32	acc_shift		Shift applied to accumulator
WORD32	input_zero_bias		Zero offset of input
WORD32	kernel_zero_bias		Zero offset of kernel
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
WORD32	out_data_format		Output data format 0:SHAPE_CUBE_DWH_T 1:SHAPE_CUBE_WHD_T
VOID *	p_scratch	xa_nn_dilated_conv2d_std_get_size()	Scratch memory pointer
WORD32	dilation_height		Kernel height dilation factor

WORD32	dilation_width		Kernel width dilation factor
Output			
WORD16 *, WORD8 *, FLOAT32 *, UWORD8 *	p_out	(out_height* out_width)* out_channels	Output cube, fixed, floating point, asym8u or asym8s, as per the out_data_format argument.

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_out, p_inp, p_ker, p_bias, p_scratch	Cannot be NULL
	Should not overlap
	Aligned on 16-byte boundary except for quantized 8-bit kernels where only p_scratch is required to be 16-byte aligned.
	For p_scratch - memory size >= size returned by xa_nn_conv2d_std_getsize()
input_height, input_width, input_channels	Greater than or equal to 1
kernel_height	{1, 2, ..., input_height}
kernel_width	{1, 2, ..., input_width}
out_channels	Greater than or equal to 1
x_stride	Equal to 1
y_stride	Equal to 1
x_padding, y_padding	Greater than or equal to 0
dilation_height, dilation_width	Greater than or equal to 1
out_height, out_width	Greater than or equal to 1
acc_shift, bias_shift, out_shift	{-31 31} for fixed point and quantized 8-bit APIs
input_zero_bias	{-255,, 0} for asym8u input, {-127,, 128} for asym8s input
kernel_zero_bias	{-255,, 0} for asym8u kernel
out_zero_bias	{0,, 255} for asym8u output, {-128,, 127} for asym8s output
out_multiplier	Greater than 0
out_data_format	Can be 0: SHAPE_CUBE_DWH_T or 1: SHAPE_CUBE_WHD_T

3.2.3 Standard 1D Convolution Kernel

Description

These kernels perform the 1D convolution operation as $z = \text{inp}(\ast)\text{kernel} + \text{bias}$. A 3D input cube (`input_height` x `input_width` x `input_channels`) is convolved with a 3D kernel cube (`kernel_height` x `input_width` x `input_channels`) to produce a 1D convolution output vector (`out_height`). With `out_channels` number of such 3D kernels, output matrix (`out_height` x `out_channels`) is produced. The bias having dimension (`out_channels`) is added after the convolution (one bias value is added to each output column) to produce the final output.

Note: The depth or channels dimension (`input_channels`) of input and kernel must be identical, and width dimension (`input_width`) of input and kernel also must be identical for 1D convolution.

`bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output, respectively. Both `bias_shift` and `acc_shift` can be either positive or negative, where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as convolution - accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note: `acc_shift` and `bias_shift` are not relevant in case of floating point kernels.

The `y_stride` argument to kernel API defines the step size of the kernel when traversing the input in height dimension.

The `y_padding` argument defines padding to the top of the input in the height dimension.

The bottom padding is calculated based on `out_height` as `bottom_paddding = kernel_height + (out_height - 1) * y_stride - (y_padding + input_height)`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

The kernel is expected to be padded if the product `input_channels*input_width` is not a multiple of 4 in case of fixed point variants, and 2 in case of floating point variant.

These kernels require temporary buffer for convolution computation. This temporary buffer is provided by `p_scratch` argument of kernel API. The size of temporary buffer should be queried using `xa_nn_conv1d_std_getsize()` helper API.

These kernels expect input and kernel cubes in `SHAPE_CUBE_DWH_T` shape type and can produce output matrix with either (`out_height` x `out_channels`) or (`out_channels` x `out_height`) dimensions. The `out_data_format` argument to kernel API controls the output matrix height and width order.

Function variants available are `xa_nn_conv1d_std_[p]`, where:

- `[p]`: precision in bits

Precision

There are four variants available:

Type	Description
16x16	16-bit kernel, 16-bit input, 16-bit output
8x16	8-bit kernel, 16-bit input, 16-bit output
8x8	8-bit kernel, 8-bit input, 8-bit output
f32	float32 kernel, float32 input, float32 output

Algorithm

$$z_{h,d} = 2^{acc-shift} \left(\sum_{i=0}^{K_H-1} \sum_{j=0}^{I_W-1} \sum_{k=0}^{I_C-1} in_{pad(h*y-stride+i),j,k} \cdot ker_{pad,d,i,j,k} + 2^{bias-shift} b_d \right)$$

$$h = 0, \dots, \overline{out-height - 1}, d = 0, \dots, \overline{out-channels - 1}$$

In case of floating point kernel, `acc_shift=0` and `bias_shift=0`.

Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

n_{pad}, ker_{pad} denote the padded `p_inp` and padded `p_ker` shapes, respectively.

K_H, I_W, I_C denote `kernel_height`, `input_width`, and `input_channels`, respectively.

b denotes the `bias` shape.

Prototype

```
WORD32 xa_nn_convld_std_getsize
(WORD32 kernel_height, WORD32 input_width, WORD32 input_channels,
 WORD32 input_precision);

WORD32 xa_nn_convld_std_16x16
(WORD16 * p_out, WORD16 * p_inp, WORD16 * p_ker,
 WORD16 * p_bias, WORD32 input_height, WORD32 input_width,
 WORD32 input_channels, WORD32 kernel_height, WORD32 out_channels,
 WORD32 y_stride, WORD32 y_padding, WORD32 out_height,
 WORD32 bias_shift, WORD32 acc_shift, WORD32 out_data_format,
 VOID * p_scratch);

WORD32 xa_nn_convld_std_8x16
(WORD16 * p_out, WORD16 * p_inp, WORD8 * p_ker,
 WORD16 * p_bias, WORD32 input_height, WORD32 input_width,
 WORD32 input_channels, WORD32 kernel_height, WORD32 out_channels,
 WORD32 y_stride, WORD32 y_padding, WORD32 out_height,
 WORD32 bias_shift, WORD32 acc_shift, WORD32 out_data_format,
 VOID * p_scratch);

WORD32 xa_nn_convld_std_8x8
(WORD8 * p_out, WORD8 * p_inp, WORD8 * p_ker,
```

```

WORD8  * p_bias,          WORD32 input_height,  WORD32 input_width,
WORD32 input_channels,    WORD32 kernel_height, WORD32 out_channels,
WORD32 y_stride,         WORD32 y_padding,     WORD32 out_height,
WORD32 bias_shift,       WORD32 acc_shift,      WORD32 out_data_format,
VOID   * p_scratch);
WORD32 xa_nn_conv1d_std_f32
(FLOAT32 * p_out,         FLOAT32 * p_inp,      FLOAT32 * p_ker,
 FLOAT32 * p_bias,        WORD32 input_height,  WORD32 input_width,
 WORD32 input_channels,   WORD32 kernel_height, WORD32 out_channels,
 WORD32 y_stride,         WORD32 y_padding,     WORD32 out_height,
 WORD32 out_data_format,  VOID   * p_scratch);

```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, FLOAT32 *,	p_inp	input_height* input_width* input_channels	Input cube, fixed or floating point, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, FLOAT32 *,	p_ker	out_channels* (kernel_height* input_width* input_channels)	Kernel cube, fixed or floating point, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, FLOAT32 *,	p_bias	out_channels	Bias vector, fixed or floating point
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Number of input channels
WORD32	kernel_height		Kernel height
WORD32	out_channels		Number of output channels
WORD32	y_stride		Vertical stride over input
WORD32	y_padding		Top padding height on input
WORD32	out_height		Output height
WORD32	bias_shift		Shift applied to bias
WORD32	acc_shift		Shift applied to accumulator
WORD32	out_data_format		Output matrix order 0: out_height x out_channels 1: out_channels x out_height
VOID *	p_scratch	xa_nn_conv1d_std_getsize()	Scratch memory pointer
Output			
WORD16 *, WORD8 *, FLOAT32 *,	p_out	out_height* out_channels	Output matrix, fixed or floating point, as per the out_data_format argument.

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_out, p_inp, p_ker, p_bias, p_scratch	Cannot be NULL
	Should not overlap
	Aligned on 16-byte boundary
	For p_scratch - memory size \geq size returned by <code>xa_nn_conv1d_std_getsize()</code>
input_height, input_width, input_channels	Greater than or equal to 1
kernel_height	{1, 2, ..., input_height}
out_channels	Greater than or equal to 1
y_stride	{1, 2, ..., kernel_height}
y_padding	Greater than or equal to 0
out_height	Greater than or equal to 1
acc_shift, bias_shift	{-31 31} for fixed point APIs
out_data_format	Can be 0: out_height x out_channels or 1: out_channels x out_height

3.2.4 Depthwise Separable 2D Convolution Kernel

Depthwise Separable 2D Convolution is computed in two steps using following two low level kernels:

- First step: `xa_nn_conv2d_depthwise_xx()` low level kernel

These kernels convolve each input 2D plane (`input_height` x `input_width`) from input cube (`input_height` x `input_width` x `input_channels`) with `channels_multiplier` number of 2D kernels (`kernel_height` x `kernel_width`) to produce `channels_multiplier` number of 2D output planes (`out_height` x `out_width`). Thus, with kernel cube of dimension (`kernel_height` x `kernel_width` x (`channels_multiplier` * `input_channels`)), output cube of dimension (`out_height` x `out_width` x (`channels_multiplier` * `input_channels`)) is produced. Bias is added to the convolution output. There is one bias value for each output 2D plane; that is, bias is a vector of dimension (`channels_multiplier` * `input_channels`).

- Second step: `xa_nn_conv2d_pointwise_xx()` low level kernel

These kernels take output cube (`out_height` x `out_width` x (`channels_multiplier` * `input_channels`)) of first step as input and perform pointwise multiplication with kernel vector (`channels_multiplier` * `input_channels`) in depth dimension to produce output 2D plane (`out_height` x `out_width`). Thus, with `out_channels` kernel vectors, output cube of dimension (`out_height` x `out_width` x `out_channels`) is produced. Bias is added to the pointwise multiplication output. There is one bias value for each output 2D plane; that is, bias is a vector of dimension `out_channels`.

Following are the descriptions for these two low level kernels.

3.2.4.1 Depthwise 2D Convolution Kernel

Description

These kernels perform the 2D depthwise convolution operation as $z = \text{inp} (*) \text{kernel} + \text{bias}$. These kernels convolve each input 2D plane ($\text{input_height} \times \text{input_width}$) from input cube ($\text{input_height} \times \text{input_width} \times \text{input_channels}$) with $\text{channels_multiplier}$ number of 2D kernels ($\text{kernel_height} \times \text{kernel_width}$) to produce $\text{channels_multiplier}$ number of 2D output planes ($\text{out_height} \times \text{out_width}$). Thus, with kernel cube of dimension ($\text{kernel_height} \times \text{kernel_width} \times (\text{channels_multiplier} * \text{input_channels})$), output cube of dimension ($\text{out_height} \times \text{out_width} \times (\text{channels_multiplier} * \text{input_channels})$) is produced. Bias is added to the convolution output. There is one bias value for each output 2D plane; that is, bias is a vector of dimension ($\text{channels_multiplier} * \text{input_channels}$).

`bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output respectively. Both `bias_shift` and `acc_shift` can be either positive or negative where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as convolution - accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note: `acc_shift` and `bias_shift` are not relevant in case of floating point kernels and quantized 8-bit kernels.

The `x_stride` and `y_stride` arguments in kernel API define the step size of the kernel when traversing the input in width and height dimensions, respectively.

The `x_padding` argument defines padding to the left of the input in the width dimension, and `y_padding` argument defines padding to the top of the input in the height dimension.

The right padding is calculated based on `out_width` as `right_paddding = kernel_width + (out_width - 1) * x_stride - (x_padding + input_width)`.

The bottom padding is calculated based on `out_height` as `bottom_paddding = kernel_height + (out_height - 1) * y_stride - (y_padding + input_height)`.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

These kernels require a temporary buffer for convolution computation. This temporary buffer is provided by the `p_scratch` argument of kernel API. The size of temporary buffer should be queried using `xa_nn_conv2d_depthwise_getsize()` helper API.

The arguments `input_zero_bias`, `kernel_zero_bias` are provided to convert the quantized 8-bit inputs into their real values and perform Depthwise 2D Convolution operation. The `out_zero_bias`, `out_multiplier` and `out_shift` values are used to quantize real values of output back to 8-bit.

The depthwise kernels expect input cube in SHAPE_CUBE_DWH_T and SHAPE_CUBE_WHD_T shape type and produce output cube in SHAPE_CUBE_DWH_T shape type respectively. The `inp_data_format` argument to the kernel API can be 0 or 1 to indicate input cube shape respectively.

The `out_data_format` argument to the kernel API must be 0 for all the kernels to indicate output cube shape.

Function variants available are `xa_nn_conv2d_depthwise_[p]`, where:

- `[p]`: precision in bits

Precision

There are six variants available:

Type	Description
16x16	16-bit kernel, 16-bit input, 16-bit output
8x16	8-bit kernel, 16-bit input, 16-bit output
8x8	8-bit kernel, 8-bit input, 8-bit output
f32	float32 kernel, float32 input, float32 output
asym8u	asym8u kernel, asym8u input, asym8u output
per_chan_sym8sxasym8s	per channel quantized sym8s kernel, asym8s input, asym8s output

Algorithm

$$z_{h,w,d \cdot C_M+m} = 2^{acc-shift} \left(\sum_{i=0}^{K_H-1} \sum_{j=0}^{K_W-1} in_{pad(h \cdot y-stride+i, (w \cdot x-stride+j), d)} \cdot ker_{pad(i, j, (d \cdot C_M+m))} + 2^{bias-shift} b_{0,0,d \cdot C_M+m} \right)$$

$$h = 0, \dots, \overline{out-height - 1}, w = 0, \dots, \overline{out-width - 1},$$

$$d = 0, \dots, \overline{input-channels - 1},$$

$$m = 0, \dots, \overline{channels-multiplier - 1}$$

In case of floating-point kernel and quantized 8-bit kernels, `acc_shift=0` and `bias_shift=0`.

Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

in_{pad}, ker_{pad} denote the padded `p_inp` and padded `p_ker` shapes, respectively.

K_H, K_W, C_M denote `kernel_height`, `kernel_width`, and `channels_multiplier`, respectively.

b denotes the `bias` shape.

Prototype

```
WORD32 xa_nn_conv2d_depthwise_getsize
(WORD32 input_height,      WORD32 input_width,      WORD32 input_channels,
 WORD32 kernel_height,    WORD32 kernel_width,    WORD32 channels_multiplier,
```

```

WORD32 x_stride,          WORD32 y_stride,          WORD32 x_padding,
WORD32 y_padding,         WORD32 output_height,     WORD32 output_width,
WORD32 circ_buf_precision, WORD32 inp_data_format);

WORD32 xa_nn_conv2d_depthwise_16x16
(WORD16 * p_out,          WORD16 * p_ker,          WORD16 * p_inp,
WORD16 * p_bias,         WORD32 input_height,     WORD32 input_width,
WORD32 input_channels,   WORD32 kernel_height,   WORD32 kernel_width,
WORD32 channels_multiplier, WORD32 x_stride,       WORD32 y_stride,
WORD32 x_padding,        WORD32 y_padding,        WORD32 out_height,
WORD32 out_width,        WORD32 acc_shift,         WORD32 bias_shift,
WORD32 inp_data_format,  WORD32 out_data_format, VOID * p_scratch);

WORD32 xa_nn_conv2d_depthwise_8x16
(WORD16 * p_out,          WORD8 * p_ker,          WORD16 * p_inp,
WORD16 * p_bias,         WORD32 input_height,     WORD32 input_width,
WORD32 input_channels,   WORD32 kernel_height,   WORD32 kernel_width,
WORD32 channels_multiplier, WORD32 x_stride,       WORD32 y_stride,
WORD32 x_padding,        WORD32 y_padding,        WORD32 out_height,
WORD32 out_width,        WORD32 acc_shift,         WORD32 bias_shift,
WORD32 inp_data_format,  WORD32 out_data_format, VOID * p_scratch);

WORD32 xa_nn_conv2d_depthwise_8x8
(WORD8 * p_out,          WORD8 * p_ker,          WORD8 * p_inp,
WORD8 * p_bias,         WORD32 input_height,     WORD32 input_width,
WORD32 input_channels,   WORD32 kernel_height,   WORD32 kernel_width,
WORD32 channels_multiplier, WORD32 x_stride,       WORD32 y_stride,
WORD32 x_padding,        WORD32 y_padding,        WORD32 out_height,
WORD32 out_width,        WORD32 acc_shift,         WORD32 bias_shift,
WORD32 inp_data_format,  WORD32 out_data_format, VOID * p_scratch);

WORD32 xa_nn_conv2d_depthwise_f32
(FLOAT32 * p_out,        const FLOAT32 * p_ker,   const FLOAT32 * p_inp,
const FLOAT32 * p_bias,  WORD32 input_height,   WORD32 input_width,
WORD32 input_channels,   WORD32 kernel_height,   WORD32 kernel_width,
WORD32 channels_multiplier, WORD32 x_stride,       WORD32 y_stride,
WORD32 x_padding,        WORD32 y_padding,        WORD32 out_height,
WORD32 out_width,        WORD32 inp_data_format,   WORD32 out_data_format,
VOID * p_scratch);

WORD32 xa_nn_conv2d_depthwise_asym8uxasym8u
(pUWORD8 p_out,          const UWORD8 * p_kernel, const UWORD8 * p_inp,
const WORD32 * p_bias,   WORD32 input_height,   WORD32 input_width,
WORD32 input_channels,   WORD32 kernel_height,   WORD32 kernel_width,
WORD32 channels_multiplier, WORD32 x_stride,       WORD32 y_stride,
WORD32 x_padding,        WORD32 y_padding,        WORD32 out_height,
WORD32 out_width,        WORD32 input_zero_bias,   WORD32 kernel_zero_bias,
WORD32 out_multiplier,   WORD32 out_shift,        WORD32 out_zero_bias,
WORD32 inp_data_format,  WORD32 out_data_format, pVOID p_scratch);

WORD32 xa_nn_conv2d_depthwise_per_chan_sym8sxasym8s
(pWORD8 p_out,          const WORD8 * p_kernel, const WORD8 * p_inp,
const WORD32 * p_bias,  WORD32 input_height,   WORD32 input_width,
WORD32 input_channels,   WORD32 kernel_height,   WORD32 kernel_width,
WORD32 channels_multiplier, WORD32 x_stride,       WORD32 y_stride,
WORD32 x_padding,        WORD32 y_padding,        WORD32 out_height,
WORD32 out_width,        WORD32 input_zero_bias,   const WORD32 * p_out_multiplier,
const WORD32 * p_out_shift, WORD32 out_zero_bias,   WORD32 inp_data_format,
WORD32 out_data_format,  pVOID p_scratch);

```

Arguments

Type	Name	Size	Description
Input			

Type	Name	Size	Description
WORD16 *, WORD8 *, FLOAT32 *, const UWORD8 *, const WORD8 *	p_ker	kernel_height* kernel_width* input_channels* channels_multiplier	Kernel cube, fixed, floating point, asym8u or sym8s, in SHAPE_CUBE_D WH or SHAPE_CUBE_W HD_T
WORD16 *, WORD8 *, FLOAT32 *, const UWORD8 *, const WORD8 *	p_inp	input_height* input_width* input_channels	Input cube, fixed, floating point, asym8u or asym8s in SHAPE_CUBE_D WH or SHAPE_CUBE_W HD_T
WORD16 *, WORD8 *, FLOAT32 *, const WORD32 *	p_bias	input_channels*chann els_multiplier	Bias vector, fixed or floating point
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Number of input channels
WORD32	kernel_height		Kernel height
WORD32	kernel_width		Kernel width
WORD32	channels_multipl ier		Multiplier value for each input channel
WORD32	x_stride		Horizontal stride over input
WORD32	y_stride		Vertical stride over input
WORD32	x_padding		Left padding width on input
WORD32	y_padding		Right padding height on input
WORD32	out_height		Output height
WORD32	out_width		Output width
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	input_zero_bias		Zero offset of input
WORD32	kernel_zero_bias		Zero offset of kernel
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32 *	p_out_multiplier	input_channels*chann els_multiplier	Array of multiplier values of output

Type	Name	Size	Description
WORD32 *	p_out_shift	input_channels*channels_multiplier	Array of shift values of output
WORD32	out_zero_bias		Zero offset of output
WORD32	inp_data_format		Input and Kernel data format 0:SHAPE_CUBE_DWH_T 1:SHAPE_CUBE_WHD_T
WORD32	out_data_format		Output data format 0:SHAPE_CUBE_DWH_T
VOID *	p_scratch	xa_nn_conv2d_depthwise_getsize()	Scratch memory pointer
Output			
WORD16 *, WORD8 *, UWORD8 *, FLOAT32 *,	p_out	out_height* out width* input_channels* channels_multiplier	Output cube, fixed, floating point, asym8u or asym8s, in SHAPE_CUBE_DWH_T

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_out, p_ker, p_inp, p_bias,	Cannot be NULL
	Should not overlap
	Aligned on <size of one element> boundary
p_scratch	Cannot be NULL
	Should not overlap with other buffers
	Aligned on 16-byte boundary
	For p_scratch - memory size >= size returned by xa_nn_conv2d_depthwise_getsize()
p_out_multiplier	Cannot be NULL
	Should not overlap
	Aligned on 4-byte boundary
p_out_shift	Cannot be NULL
	Should not overlap
	Aligned on 4-byte boundary
	Each 32-bit value should be in range [-31 ... 31]

<code>input_height, input_width, input_channels</code>	Greater than or equal to 1
<code>kernel_height</code>	{1,2,...,input_height}
<code>kernel_width</code>	{1,2,...,input_width}
<code>channels_multiplier</code>	Greater than or equal to 1
<code>x_stride</code>	{1,2,...,kernel_width}
<code>y_stride</code>	{1,2,...,kernel_height}
<code>x_padding, y_padding</code>	Greater than or equal to 0
<code>out_height, out_width</code>	Greater than or equal to 1
<code>acc_shift, bias_shift, out_shift</code>	{-31 31} for fixed point and quantized 8-bit APIs
<code>input_zero_bias</code>	{-255,....., 0} for asym8u input, {-127....., 128} for asym8s input
<code>kernel_zero_bias</code>	{-255,....., 0} for asym8u kernel
<code>out_zero_bias</code>	{0,.....,255} for asym8u output, {-128....., 127} for asym8s output
<code>out_multiplier</code>	Greater than 0
<code>inp_data_format</code>	can be 0: SHAPE_CUBE_DWH_T or 1: SHAPE_CUBE_WHD_T
<code>out_data_format</code>	must be 0: SHAPE_CUBE_DWH_T

3.2.4.2 Pointwise 2D Convolution Kernel

Description

These kernels perform pointwise multiplication of input cube (`input_height` x `input_width` x `input_channels`) with kernel vector (`input_channels`) in depth dimension to produce output 2D plane (`input_height` x `input_width`). Thus, with `out_channels` kernel vectors, output cube of dimension (`input_height` x `input_width` x `out_channels`) is produced. Bias is added to the pointwise multiplication output. There is one bias value for each output 2D plane; that is, bias is a vector of dimension `out_channels`.

The `bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output respectively. Both `bias_shift` and `acc_shift` can be either positive or negative, where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as convolution - accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note: `acc_shift` and `bias_shift` are not relevant in case of floating point kernels and quantized 8-bit kernels.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

These kernels expect input cube in SHAPE_CUBE_DWH_T shape type, kernel as matrix, bias as vector and produce output cube in SHAPE_CUBE_WHD_T or SHAPE_CUBE_DWH_T (only for 8x8, asym8uxasym8u and per_chan_sym8sxsasym8s kernels) shape type. The `out_data_format` argument

to kernel API must be always 1 except for 8x8 and quantized 8-bit kernels for which it can be 0 or 1 indicating SHAPE_CUBE_DWH_T and SHAPE_CUBE_WHD_T respectively.

Function variants available are `xa_nn_conv2d_pointwise_[p]`, where:

- `[p]`: precision in bits

Precision

There are six variants available:

Type	Description
16x16	16-bit kernel, 16-bit input, 16-bit output
8x16	8-bit kernel, 16-bit input, 16-bit output
8x8	8-bit kernel, 8-bit input, 8-bit output
f32	float32 kernel, float32 input, float32 output
asym8u/asym8s	asym8u kernel, asym8u input, asym8u output
per_chan_sym8s/asym8s	per channel quantized sym8s kernel, asym8s input, asym8s output

Algorithm

$$z_{h,w,d} = 2^{acc-shift} \left(\sum_{k=0}^{I_C-1} in_{h,w,k} \cdot ker_{d,0,0,k} + 2^{bias-shift} b_{0,0,d} \right)$$

$$h = 0, \dots, \overline{input-height} - 1, w = 0, \dots, \overline{input-width} - 1,$$

$$d = 0, \dots, \overline{out_channels} - 1$$

In case of floating-point kernel and quantized 8-bit kernels, `acc_shift=0` and `bias_shift=0`. Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

`in`, `ker` denote the `p_inp`, and `p_ker` shapes respectively.

`IC` denotes `input_channels`

`b` denotes the `bias` shape

Prototype

```
WORD32 xa_nn_conv2d_pointwise_16x16
(WORD16 * p_out,          WORD16 * p_ker,          WORD16 * p_inp,
 WORD16 * p_bias,         WORD32 input_height,     WORD32 input_width,
 WORD32 input_channels,   WORD32 out_channels,     WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 out_data_format);
WORD32 xa_nn_conv2d_pointwise_8x16
(WORD16 * p_out,          WORD8 * p_ker,           WORD16 * p_inp,
 WORD16 * p_bias,         WORD32 input_height,     WORD32 input_width,
 WORD32 input_channels,   WORD32 out_channels,     WORD32 acc_shift,
 WORD32 bias_shift,       WORD32 out_data_format);
WORD32 xa_nn_conv2d_pointwise_8x8
(WORD8 * p_out,           WORD8 * p_ker,           WORD8 * p_inp,
 WORD8 * p_bias,          WORD32 input_height,     WORD32 input_width,
```

```

WORD32 input_channels,      WORD32 out_channels,      WORD32 acc_shift,
WORD32 bias_shift,         WORD32 out_data_format);
WORD32 xa_nn_conv2d_pointwise_f32
(FLOAT32 * p_out,          const FLOAT32 * p_ker, const FLOAT32 * p_inp,
 const FLOAT32 * p_bias,   WORD32 input_height,   WORD32 input_width,
 WORD32 input_channels,    WORD32 out_channels,
 WORD32 out_data_format);
WORD32 xa_nn_conv2d_pointwise_asym8uxasym8u
(UWORD8 * p_out,          const UWORD8 * p_ker, const UWORD8 * p_inp,
 WORD32 * p_bias,         WORD32 input_height,   WORD32 input_width,
 WORD32 input_channels,    WORD32 out_channels,   WORD32 input_zero_bias,
 WORD32 kernel_zero_bias, WORD32 out_multiplier, WORD32 out_shift,
 WORD32 out_zero_bias,     WORD32 out_data_format);
WORD32 xa_nn_conv2d_pointwise_asym8uxasym8u
(UWORD8 * p_out,          const UWORD8 * p_ker, const UWORD8 * p_inp,
 const WORD32 * p_bias,   WORD32 input_height,   WORD32 input_width,
 WORD32 input_channels,    WORD32 out_channels,   WORD32 input_zero_bias,
 WORD32 kernel_zero_bias, WORD32 out_multiplier, WORD32 out_shift,
 WORD32 out_zero_bias,     WORD32 out_data_format);
WORD32 xa_nn_conv2d_pointwise_per_chan_sym8sxsasym8s
(WORD8 * p_out,          const WORD8 * p_ker, const WORD8 * p_inp,
 const WORD32 * p_bias,   WORD32 input_height,   WORD32 input_width,
 WORD32 input_channels,    WORD32 out_channels,   WORD32 input_zero_bias,
 WORD32 * p_out_multiplier, WORD32 * p_out_shift, WORD32 out_zero_bias,
 WORD32 out_data_format);

```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_ker	out_channels * input_channels	Kernel matrix, fixed, floating point, asym8u or asym8s, (out_channels x input_channels)
WORD16 *, WORD8 *, const FLOAT32 *, const UWORD8 *, const WORD8 *	p_inp	input_height* input width* input_channels	Input cube, fixed or floating point, asym8u or sym8s, in SHAPE_CUBE_DWH_T
WORD16 *, WORD8 *, FLOAT32 *, const WORD32 *	p_bias	out_channels	Bias vector, fixed or floating point
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Number of input channels
WORD32	out_channels		Number of output channels
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	input_zero_bias		Zero offset of input

WORD32	kernel_zero_bias		Zero offset of kernel
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
WORD32	out_data_format		Output data format 0:SHAPE_CUBE_DWH_T 1:SHAPE_CUBE_WHD_T
Output			
WORD16 *, WORD8 *, FLOAT32 *, UWORD8 *	p_out	(out_height* out_width)* out_channels	Output cube, fixed, floating point, asym8u or asym8s, as per the out_data_format argument.

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_out, p_ker, p_inp, p_bias	Cannot be NULL
	Should not overlap
	Aligned on 16-byte boundary except for 8x8 and quantized 8-bit kernels
input_height, input_width	Greater than or equal to 1
input_channels	Greater than or equal to 4, multiple of 4 except for 8x8 and asym8u kernels
out_channels	Greater than or equal to 1
acc_shift, bias_shift, out_shift	{-31 31} for fixed point and quantized 8-bit APIs
input_zero_bias	{-255,....., 0} for asym8u input, {-127....., 128} for asym8s input
kernel_zero_bias	{-255....., 0} for asym8u kernel
out_zero_bias	{0,.....,255} for asym8u output, {-128....., 127} for asym8s output
out_multiplier	Greater than 0
out_data_format	Can be 0: SHAPE_CUBE_DWH_T or 1: SHAPE_CUBE_WHD_T for 8x8 and quantized 8-bit kernels. Must be 1 for other kernels.

3.3 Activation Kernels

3.3.1 Sigmoid

Description

These kernels perform the sigmoid operation on input vector x and give output vector as $y = \text{sigmoid}(x)$. Both the input and output vectors have size `vec_length`.

The fixed-point kernels accept 32-bit input in Q6.25 format and give output in Q16.15 (32-bit), Q15 (16-bit), or Q7 (8-bit) format.

For the `asym8u` and `asym8s` kernels both the input and output are of `asym8u` and `asym8s` datatype respectively.

Function variants available are `xa_nn_vec_sigmoid_[p]_[q]`, where:

- `[p]`: Input precision in bits
- `[q]`: Output precision in bits

Precision

There are six variants available.

Type	Description
<code>32_32</code>	32-bit input, 32-bit output
<code>32_16</code>	32-bit input, 16-bit output
<code>32_8</code>	32-bit input, 8-bit output
<code>f32_f32</code>	float32 input, float32 output
<code>asym8u_x_asym8u</code>	<code>asym8u</code> input, <code>asym8u</code> output
<code>asym8s_x_asym8s</code>	<code>asym8s</code> input, <code>asym8s</code> output

Algorithm

$$y_n = \frac{1}{1 + \exp(-x_n)}, \quad n = 0, \dots, \overline{vec_length - 1}$$

Prototype

```
WORD32 xa_nn_vec_sigmoid_32_32
(WORD32 * p_out,          const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_sigmoid_32_16
(WORD16 * p_out,          const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_sigmoid_32_8
(WORD8 * p_out,           const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_sigmoid_f32_f32
(FLOAT32 * p_out,         const FLOAT32 * p_vec,     WORD32 vec_length);
WORD32 xa_nn_vec_sigmoid_asym8u_asym8u
(UWORD8 * p_out,          const UWORD8 * p_vec,      WORD32 zero_point,
WORD32 input_range_radius, WORD32 input_multiplier, WORD32 input_left_shift,
```

```
WORD32 vec_length);
WORD32 xa_nn_vec_sigmoid_asym8s_asym8s
(WORD8 * p_out, const WORD8 * p_vec, WORD32 zero_point,
WORD32 input_range_radius, WORD32 input_multiplier, WORD32 input_left_shift,
WORD32 vec_length);
```

Arguments

Type	Name	Size	Description
Input			
const WORD32 *, const UWORD8 *, const FLOAT32 *, const WORD8 *	p_vec	vec_length	Input vector, Q6.25, floating point, asym8u or asym8s
WORD32	zero_point		bias value
WORD32	input_range_radius		Range radius: For asym8u output = ((x _i - zero_point) < radius)? sigmoid() : 255 output = ((x _i - zero_point) > (-radius))? sigmoid() : 0 For asym8s output = ((x _i - zero_point) < radius)? sigmoid() : 127 output = ((x _i - zero_point) > (-radius))? sigmoid() : -128
WORD32	input_multiplier		Multiplier value of input
WORD32	input_left_shift		Left Shift value of input
WORD32	vec_length		Length of input vector
Output			
WORD32 *, WORD16 *, WORD8 *, UWORD8 *, FLOAT32 *	p_out	vec_length	Output vector, fixed (Q16.15, Q15, Q7), floating point, asym8u or asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_vec, p_out	Should not overlap Cannot be NULL
zero_point	[0, 255] for asym8u [-128, 127] for asym8s
input_range_radius	[0, 255]
input_left_shift	[-31, 31]
input_multiplier	Shouldn't be less than 0.

vec_length	Greater than 0
------------	----------------

3.3.2Tanh

Description

These kernels perform the hyperbolic tangent operation on input vector x and give output vector as $y = \tanh(x)$. Both the input and output vectors have size `vec_length`.

The fixed-point kernels accept 32-bit input in Q6.25 format and give output in Q16.15 (32-bit), Q15 (16-bit), or Q7 (8-bit) format.

For the `asym8s` kernels both the input and output are of `asym8s` datatype.

Function variants available are `xa_nn_vec_tanh_[p]_[q]`, where:

- `[p]`: Input precision in bits
- `[q]`: Output precision in bits

Precision

There are five variants available:

Type	Description
<code>32_32</code>	32-bit input, 32-bit output
<code>32_16</code>	32-bit input, 16-bit output
<code>32_8</code>	32-bit input, 8-bit output
<code>f32_f32</code>	float32 input, float32 output
<code>asym8sxasym8s</code>	asym8s input, asym8s output

Algorithm

$$y_n = \tanh(x_n), \quad n = 0, \dots, \overline{vec_length} - 1$$

Prototype

```
WORD32 xa_nn_vec_tanh_32_32
(WORD32 * p_out,          const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_tanh_32_16
(WORD16 * p_out,          const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_tanh_32_8
(WORD8 * p_out,           const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_tanh_f32_f32
(FLOAT32 * p_out,         const FLOAT32 * p_vec,    WORD32 vec_length);
WORD32 xa_nn_vec_tanh_asym8s_asym8s
(WORD8 * p_out,           const WORD8 * p_vec,      WORD32 zero_point,
WORD32 input_range_radius, WORD32 input_multiplier, WORD32 input_left_shift,
WORD32 vec_length);
```

Arguments

Type	Name	Size	Description
Input			
const WORD32 *, const FLOAT32 *, const WORD8 *	p_vec	vec_length	Input vector, Q6.25, floating point or asym8s
WORD32	zero_point		Bias value
WORD32	input_range_radius		Range radius: output = ((x _i - zero_point) < radius)? tanh() : 127 output = ((x _i - zero_point) > (-radius))? tanh() : -128
WORD32	input_multiplier		Multiplier value of input
WORD32	input_left_shift		Left shift value of input
WORD32	vec_length		Length of input vector
Output			
WORD32 *, WORD16 *, WORD8 *, FLOAT32 *	p_out	vec_length	Output vector, fixed (Q16.15, Q15, Q7), floating point or asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_vec, p_out	Should not overlap
	Cannot be NULL
zero_point	[-128, 127]
input_range_radius	[0, 255]
input_multiplier	Shouldn't be less than 0
vec_length	Greater than 0

3.3.3 Rectifier Linear Unit (ReLU)

Description

These kernels compute the rectifier linear unit function of input vector x and give output vector as $y = \text{relu}(x)$. Both the input and output vectors have size `vec_length`.

The fixed-point routines accept 32-bit input in Q6.25 format and gives 32-bit output in Q16.15 format.

The `threshold` argument to `relu` kernel API allows to set upper threshold for proper compression of output signal and is expected in Q16.15 format. In `relu1` and `relu6` kernels, the thresholds are set to 1 and 6, respectively.

For the `asym8u` and `asym8s` kernels, the quantized input is requantized and applied the standard ReLU function to give the output. The `threshold` argument is not applicable for quantized ReLU kernels.

The standard ReLU kernels `relu_std` can be used when the `threshold` is not required.

Function variants available are `xa_nn_vec_relu_[p]_[q]`, `xa_nn_vec_relu1_[p]_[q]`, and `xa_nn_vec_relu6_[p]_[q]`, where:

- `[p]`: Input precision in bits
- `[q]`: Output precision in bits

Precision

There are six variants available:

Type	Description
<code>32_32</code>	32-bit input, 32-bit output
<code>f32_f32</code>	float32 input, float32 output
<code>16_16</code>	16-bit input, 16-bit output
<code>8_8</code>	8-bit input, 8-bit output
<code>asym8u_asym8u</code>	asym8u input, asym8u output
<code>asym8s_asym8s</code>	asym8s input, asym8s output

Algorithm

$$y_n = \max(0, \min(x_n, K)), \quad n = 0, \dots, \overline{vec_length} - 1$$

K represents threshold

Prototype

```
WORD32 xa_nn_vec_relu_32_32
(WORD32 * p_out,      const WORD32 * p_vec,   WORD32 threshold,
 WORD32 vec_length);
WORD32 xa_nn_vec_relu_f32_f32
(FLOAT32 * p_out,     const FLOAT32 * p_vec,   FLOAT32 threshold,
 WORD32 vec_length);
```

```

WORD32 xa_nn_vec_relu_16_16
(WORD16 * p_out,      const WORD16 * p_vec,   WORD16 threshold,
 WORD32 vec_length);
WORD32 xa_nn_vec_relu_8_8
(WORD8 * p_out,      const WORD8 * p_vec,   WORD8 threshold,
 WORD32 vec_length);
WORD32 xa_nn_vec_relu_asym8u_asym8u
(UWORD8 * p_out,      const UWORD8 * p_vec, WORD32 inp_zero_bias,
 WORD32 out_multiplier, WORD32 out_shift,   WORD32 out_zero_bias,
 WORD32 quantized_activation_min, WORD32 quantized_activation_max,
 WORD32 vec_length);
WORD32 xa_nn_vec_relu_asym8s_asym8s
(WORD8 * p_out,      const WORD8 * p_vec, WORD32 inp_zero_bias,
 WORD32 out_multiplier, WORD32 out_shift,   WORD32 out_zero_bias,
 WORD32 quantized_activation_min, WORD32 quantized_activation_max,
 WORD32 vec_length);
WORD32 xa_nn_vec_relu1_32_32
(WORD32 * p_out,      const WORD32 * p_vec,   WORD32 vec_length);
WORD32 xa_nn_vec_relu1_f32_f32
(FLOAT32 * p_out,      const FLOAT32 * p_vec, WORD32 vec_length);
WORD32 xa_nn_vec_relu6_32_32
(WORD32 * p_out,      const WORD32 * p_vec,   WORD32 vec_length);
WORD32 xa_nn_vec_relu6_f32_f32
(FLOAT32 * p_out,      const FLOAT32 * p_vec, WORD32 vec_length);
WORD32 xa_nn_vec_relu_std_32_32
(WORD32 * p_out,      const WORD32 * p_vec,   WORD32 vec_length);
WORD32 xa_nn_vec_relu_std_f32_f32
(FLOAT32 * p_out,      const FLOAT32 * p_vec, WORD32 vec_length);
WORD32 xa_nn_vec_relu_std_16_16
(WORD16 * p_out,      const WORD16 * p_vec,   WORD32 vec_length);
WORD32 xa_nn_vec_relu_std_8_8
(WORD8 * p_out,      const WORD8 * p_vec,   WORD32 vec_length);

```

Arguments

Type	Name	Size	Description
Input			
const WORD32 *, const FLOAT32 *, const WORD16 *, const WORD8 *, const UWORD8 *	p_vec	vec_length	Input vector, fixed-point, floating point, asym8u or asym8s
WORD32	inp_zero_bias		Zero bias value for input vector
WORD32	out_multiplier		Fixed-point multiplier value for output
WORD32	out_shift		Shift value for output
WORD32	vec_length		length of input vector
WORD32	out_zero_bias		Zero bias value for output vector
WORD32	quantized_act ivation_min		Lower threshold value, quantized.
WORD32, FLOAT32	quantized_act ivation_max		Upper threshold value, quantized
WORD32 FLOAT32 WORD16 WORD8	threshold		threshold, fixed or floating point

Output			
WORD32 *, FLOAT32 *, WORD16 *, WORD8 *, UWORD8 *	p_out	vec_length	Output vector, fixed-point, floating point, asym8u or asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_vec, p_out	Should not overlap Cannot be NULL
inp_zero_bias, out_zero_bias	{0,.....,255} for asym8u, {-128....., 127} for asym8s input
out_multiplier	Shouldn't be less than 0.
out_shift	{-31, ..., 31}
quantized_activation_min quantized_activation_max	{0,.....,255} for asym8u output, {-128....., 127} for asym8s output quantized_activation_min < quantized_activation_max

3.3.4 Softmax

Description

These kernels compute the softmax (normalized exponential function) of input vector x and give output vector as $y = \text{softmax}(x)$. Both the input and output vectors have size `vec_length`.

The fixed-point kernels accept 32-bit input in Q6.25 format and give 32-bit output in Q16.15 format.

For the `asym8u` kernels, both the input and output are of the same precision and for `asym8s` kernels, the input is `asym8s` and the output precision can be `asym8s` or 16-bit.

Function variants available are `xa_nn_vec_softmax_[p]_[q]`, where:

- `[p]`: Input precision in bits
- `[q]`: Output precision in bits

Precision

There are five variants available:

Type	Description
<code>32_32</code>	32-bit input, 32-bit output
<code>f32_f32</code>	float32 input, float32 output
<code>asym8u_asym8u</code>	asym8u input, asym8u output
<code>asym8s_asym8s</code>	asym8s input, asym8s output
<code>asym8s_16</code>	asym8s input, 16-bit output

Algorithm

$$y_n = \frac{\exp(x_n)}{\sum_k \exp(x_k)}, \quad n = 0, \dots, \overline{vec_length} - 1$$

Prototype

```
WORD32 xa_nn_vec_softmax_32_32
(WORD32 * p_out,      const WORD32 * p_vec,      WORD32 vec_length);
WORD32 xa_nn_vec_softmax_f32_f32
(FLOAT32 * p_out,     const FLOAT32 * p_vec,     WORD32 vec_length);
WORD32 xa_nn_vec_softmax_asym8u_asym8u
(UWORD8 * p_out,      const UWORD8 * p_vec,      WORD32 diffmin,
WORD32 input_left_shift, WORD32 input_multiplier,
WORD32 vec_length,    pVOID p_scratch);
WORD32 xa_nn_vec_softmax_asym8s_asym8s
(WORD8 * p_out,       const WORD8 * p_vec,       WORD32 diffmin,
WORD32 input_left_shift, WORD32 input_multiplier,
WORD32 vec_length,    pVOID p_scratch);
WORD32 xa_nn_vec_softmax_asym8s_16
(WORD16 * p_out,      const WORD8 * p_vec,      WORD32 diffmin,
WORD32 input_left_shift, WORD32 input_multiplier,
WORD32 vec_length,    pVOID p_scratch);
```

Arguments

Type	Name	Size	Description
Input			
WORD32 *, FLOAT32 *, const UWORD8 *, const WORD8 *	p_vec	vec_length	Input vector, Q6.25, floating point, asym8u or asym8s
WORD32	diffmin		Diffmin value: output = $((x_i - \max) > \text{diffmin}) ? \text{softmax}() : 0$
WORD32	input_left_shift		left shift value of input
WORD32	input_multiplier		multiplier value of input
WORD32	vec_length		Length of input vector
Output			
WORD32 *, FLOAT32 *, UWORD8 *, WORD8 *, WORD16 *	p_out	vec_length	Output vector, Q16.15, floating point, asym8u, asym8s or 16-bit.
Temporary			
VOID *	p_scratch		Scratch (temporary) memory pointer

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_vec, p_out	Should not overlap
	Cannot be NULL
input_left_shift	[-31, 31]
input_multiplier	Shouldn't be less than 0.
vec_length	Greater than 0

3.3.5 Activation Min Max

Description

These kernels compute the activation minimum and maximum value of input vector x and give output vector as $y = \text{activation_min_max}(x)$. Both the input and output vectors have size `num_elm`.

For activation min max kernels, the input precision and the output precision are same.

The `activation_min` and `activation_max` arguments to the kernel API allow to set the threshold for proper compression of the output. The kernel is a generic implementation of the ReLU function.

Function variant available is `xa_nn_vec_activation_min_max_[p]_[q]`, where:

- `[p]`: Input precision in bits
- `[q]`: Output precision in bits

Precision

There are four variants available:

Type	Description
<code>f32_f32</code>	float32 input, float32 output
<code>asym8u_asym8u</code>	asym8u input, asym8u output
<code>16_16</code>	16-bit input, 16-bit output
<code>8_8</code>	8-bit input, 8-bit output

Algorithm

$$y_n = \max(\text{activation_min}, \min(x_n, \text{activation_max})), \quad n = 0, \dots, \text{vec_length} - 1$$

`activation_min` represents lower threshold.

`activation_max` represents upper threshold.

Prototype

```
WORD32 xa_nn_vec_activation_min_max_f32_f32
(FLOAT32 * p_out,      const FLOAT32 * p_vec, FLOAT32 activation_min,
 FLOAT32 activation_max, WORD32  vec_length);
WORD32 xa_nn_vec_activation_min_max_asym8u_asym8u
(UWORD8 * p_out,      const UWORD8 * p_vec,  int activation_min,
 int activation_max,   WORD32  vec_length);
WORD32 xa_nn_vec_activation_min_max_16_16
(WORD16 * p_out,      const WORD16 * p_vec,  int activation_min,
 int activation_max,   WORD32  vec_length);
WORD32 xa_nn_vec_activation_min_max_8_8
(WORD8 * p_out,      const WORD8 * p_vec,   int activation_min,
 int activation_max,   WORD32  vec_length);
```

Arguments

Type	Name	Size	Description
Input			
const UWORD8 *, const FLOAT32 *, const WORD16 *, const WORD8 *	<code>p_vec</code>	<code>vec_length</code>	Input vector, floating-point, asym8u or fixed point.
WORD32	<code>vec_length</code>		Length of input vector
WORD32, FLOAT32	<code>activation_min</code>		Lower threshold value, floating-point, asym8u or fixed point.
WORD32, FLOAT32	<code>activation_max</code>		Upper threshold value, floating-point, asym8u or fixed point

Output			
UWORD8 *, FLOAT32 *, WORD16 *, WORD8 *	p_out	vec_length	Output vector, floating-point, asym8u or fixed point

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_vec, p_out	Aligned on (size of one element)-byte boundary Cannot be NULL

3.3.6 Hard Swish

Description

These kernels compute the hard-swish function of input vector x and give output vector as $y = \text{hard_swish}(x)$. Both the input and output vectors have size `vec_length`.

The hard-swish activation function is a type of activation function based on swish but replaces the computationally expensive sigmoid function by ReLU6.

Function variants available are `xa_nn_vec_hard_swish_[p]_[q]`, where:

- [p]: Input precision in bits
- [q]: Output precision in bits

Precision

There is one variant available:

Type	Description
asym8s_asym8s	asym8s input, asym8s output

Algorithm

$$y_n = x_n * [\text{ReLU6}(x_n + 3)/6], \quad n = 0, \dots, \overline{vec_length - 1}$$

Prototype

```
WORD32 xa_nn_vec_hard_swish_asym8s_asym8s
(WORD8 * p_out, const WORD8 * p_vec, WORD32 inp_zero_bias,
WORD16 reluish_multiplier, WORD32 reluish_shift, WORD16 out_multiplier,
WORD32 out_shift, WORD32 out_zero_bias, WORD32 vec_length);
```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *	p_vec	vec_length	Input vector, asym8s
WORD32	inp_zero_bias		Zero bias value for input vector
WORD16	reluish_multiplier		Fixed-point multiplier value for reluish scale
WORD32	reluish_shift		Shift value for reluish scale
WORD16	out_multiplier		Fixed-point multiplier value for output
WORD32	out_shift		Shift value for output
WORD32	out_zero_bias		Zero bias value for output vector
WORD32	vec_length		length of input vector
Output			
WORD8 *	p_out	vec_length	Output vector, asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_vec, p_out	Cannot be NULL Should not overlap (the two pointers could be same, inplace operation is possible)
inp_zero_bias, out_zero_bias	{-128....., 127} for asym8s datatype
out_multiplier, reluish_multiplier	Shouldn't be less than 0
out_shift, reluish_shift	{-31, ..., 31}

3.3.7 Parametric ReLU (PReLU)

Description

These kernels compute the Parametric ReLU function of input vector x and give output vector as $y = \text{prelu}(x)$. Both the input and output vectors have size `vec_length`.

The PReLU activation function acts like a standard ReLU function for input values greater than or equal to 0. For input values less than 0, a learnable parameter $\alpha(a)$ is multiplied with input to get the output.

Function variants available are `xa_nn_vec_prelu_[p]_[q]`, where:

- [p]: Input precision in bits
- [q]: Output precision in bits

Precision

There is one variant available:

Type	Description
asym8s_asym8s	asym8s input, asym8s output

Algorithm

$$y_n = x_n, \quad \text{when } x_n \geq 0 \quad n = 0, \dots, \text{vec_length} - 1$$

$$y_n = ax_n, \quad \text{when } x_n < 0$$

where a is the learnable parameter alpha.

Prototype

```
WORD32 xa_nn_vec_prelu_asym8s_asym8s
(WORD8 * p_out, const WORD8 * p_vec, const WORD8 * p_vec_alpha,
 WORD32 inp_zero_bias, WORD32 alpha_zero_bias, WORD32 alpha_multiplier,
 WORD32 alpha_shift, WORD32 out_multiplier, WORD32 out_shift,
 WORD32 out_zero_bias, WORD32 vec_length);
```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *	p_vec	vec_length	Input vector, asym8s
const WORD8 *	p_vec_alpha	vec_length	alpha input vector, asym8s
WORD32	inp_zero_bias		Zero bias value for input vector
WORD32	alpha_zero_bias		Zero bias value for alpha input vector
WORD16	alpha_multiplier		Fixed-point multiplier value for alpha input.
WORD32	alpha_shift		Shift value for alpha input.
WORD16	out_multiplier		Fixed-point multiplier value for output
WORD32	out_shift		Shift value for output
WORD32	out_zero_bias		Zero bias value for output vector
WORD32	vec_length		length of input vector
Output			
WORD8 *	p_out	vec_length	Output vector, asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
<code>p_vec, p_out</code>	Cannot be NULL Should not overlap (the two pointers could be same, inplace operation is possible)
<code>inp_zero_bias, alpha_zero_bias</code>	{-127....., 128} for asym8s datatype
<code>out_zero_bias</code>	{-128....., 127} for asym8s datatype
<code>out_multiplier, alpha_multiplier</code>	Shouldn't be less than 0
<code>out_shift, alpha_shift</code>	{-31, ..., 31}

3.4 Pooling Kernels

3.4.1 Average Pool Kernel

Description

These kernels compute 2D average pool on a set of input planes (matrices) x and give a set of planes y as output.

The pooling region is defined by `kernel_height` and `kernel_width`. It is shifted over the input plane in steps of `x_stride` horizontally and in steps of `y_stride` vertically to generate the specified output plane size. The input is extended by zero padding as specified by the padding region. The padding is determined by the parameters `x_padding`, `y_padding` for left and top side padding respectively, and `out_width`, `out_height` for right and bottom padding respectively. Around the edges of input planes, if only a part of pooling region is covering input plane then only the average of those elements is calculated, and the denominator is the number of elements from input in current pooling region.

The average pool kernels accept input as 8-bit, 16-bit integer, asym8u or single precision floating point format and give output in the same precision as input.

These kernels require temporary buffer for average pool computation. This temporary buffer is provided by the `p_scratch` argument of kernel API. The size of the temporary buffer should be queried using `xa_nn_avgpool_getsize()` helper API.

These kernels expect input cube in `SHAPE_CUBE_WHD_T` and `SHAPE_CUBE_DWH_T` shape type and produce output cube in `SHAPE_CUBE_WHD_T` and `SHAPE_CUBE_DWH_T` shape type respectively. The `inp_data_format` and `out_data_format` arguments to the kernel API can be 0 or 1 to indicate input and output cube shapes respectively.

The value of `inp_data_format` and `out_data_format` must be equal.

Note, the fixed-point 8-bit average pool kernel, `xa_nn_avgpool_8` can be used for the quantized int8 datatype.

Function variants available are `xa_nn_avgpool_[p]`, where:

- `[p]`: Input and Output precision in bits

Precision

There are four variants available:

Type	Description
8	8-bit input, 8-bit output
16	16-bit input, 16-bit output
f32	float32 input, float32 output
asym8u	asym8u input, asym8u output

Algorithm

$$z_{h,w,d} = \frac{1}{K_H K_W} \left(\sum_{i=0}^{K_H-1} \sum_{j=0}^{K_W-1} in_{(h*y-stride+i),(w*x-stride+j),d} \right)$$

$$h = 0, \dots, \overline{out_height - 1}, \quad w = 0, \dots, \overline{out_width - 1},$$

$$d = 0, \dots, \overline{out_channels - 1}$$

in denotes padded input cube, *z* denotes output

K_H, K_W denote kernel_height, kernel_width respectively.

Prototype

```
WORD32 xa_nn_avgpool_getsize
(WORD32 input_channels, WORD32 inp_precision, WORD32 out_precision,
 WORD32 input_height, WORD32 input_width, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 inp_data_format, WORD32 out_data_format);

WORD32 xa_nn_avgpool_8
(WORD8 * p_out, const WORD8 * p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 inp_data_format, WORD32 out_data_format,
 VOID * p_scratch);

WORD32 xa_nn_avgpool_16
(WORD16 * p_out, const WORD16 * p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 inp_data_format, WORD32 out_data_format,
 VOID * p_scratch);

WORD32 xa_nn_avgpool_f32
(FLOAT32 * p_out, const FLOAT32 * p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 inp_data_format, WORD32 out_data_format,
```



```

VOID * p_scratch);
WORD32 xa_nn_avgpool_asym8u
(UWORD8* p_out,      const UWORD8* p_inp,      WORD32 input_height,
WORD32 input_width,  WORD32 input_channels,    WORD32 kernel_height,
WORD32 kernel_width, WORD32 x_stride,          WORD32 y_stride,
WORD32 x_padding,    WORD32 y_padding,         WORD32 out_height,
WORD32 out_width,    WORD32 inp_data_format,   WORD32 out_data_format,
VOID *p_scratch);

```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *, const WORD16 *, const UWORD8 *, const FLOAT32 *	p_inp	input_height * input_width * input_channels	Input cube
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Input number of channels
WORD32	kernel_height		Pooling window height
WORD32	kernel_width		Pooling window width
WORD32	x_stride		Horizontal stride over input
WORD32	y_stride		Vertical stride over input
WORD32	x_padding		Left padding width on input
WORD32	y_padding		Top padding height on input
WORD32	out_height		Output height
WORD32	out_width		Output width
WORD32	inp_data_format		Input data format 0: SHAPE_CUBE_DWH_T 1: SHAPE_CUBE_WHD_T
WORD32	out_data_format		Output data format: 0: SHAPE_CUBE_DWH_T 1: SHAPE_CUBE_WHD_T
Output			
WORD8 *, WORD16 *, UWORD8 *, FLOAT32 *	p_out	out_height * out_width * input_channels	Output cube
Temporary			
VOID *	p_scratch	xa_nn_avgpool_ getsize()	Temporary / scratch memory

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
<code>p_inp, p_out</code>	Cannot be NULL Should not overlap
<code>p_scratch</code>	Cannot be NULL Should not overlap Memory size \geq size returned by <code>xa_nn_avgpool_getsize()</code>
<code>input_height, input_width</code>	Greater than or equal to 1
<code>input_channels</code>	Greater than or equal to 1
<code>kernel_height</code>	{1, 2, ..., min(input_height, 256)} (for 8-bit and 16-bit) {1, 2, ..., input_height} (for float32)
<code>kernel_width</code>	{1, 2, ..., min(input_width, 256)} (for 8-bit and 16-bit) {1, 2, ..., input_width} (for float32)
<code>x_stride, y_stride</code>	Greater than or equal to 1
<code>x_padding, y_padding</code>	Greater than or equal to 0
<code>out_height, out_width</code>	greater than or equal to 1
<code>inp_data_format</code>	Can be 0: SHAPE_CUBE_DWH_T or 1: SHAPE_CUBE_WHD_T
<code>out_data_format</code>	Must be equal to <code>inp_data_format</code>

3.4.2Max Pool Kernel

Description

These kernels perform 2D max pooling operation over a set of input planes x and give as output, a set of planes y .

The pooling region is defined by `kernel_height` and `kernel_width`. It is shifted over the input plane horizontally in steps of `x_stride` and vertically in steps of `y_stride` to generate the specified output plane size.

The input plane, padded with the maximum negative values is considered while performing the max pooling operation. The padding region is determined by the parameters `x_padding`, `y_padding` for left and top side padding respectively, and `out_width`, `out_height` for right and bottom padding respectively.

The max pool kernels accept input as 8-bit, 16-bit integer, `asym8u` or single precision floating point format and give output in the same precision as input.

These kernels require temporary buffer for max pool computation. This temporary buffer is provided by the `p_scratch` argument of kernel API. The size of the temporary buffer should be queried using the `xa_nn_maxpool_getsize()` helper API.

These kernels expect input cube in SHAPE_CUBE_WHD_T and SHAPE_CUBE_DWH_T shape type and produce output cube in SHAPE_CUBE_WHD_T and SHAPE_CUBE_DWH_T shape type respectively. The

`inp_data_format` and `out_data_format` arguments to the kernel API can be 0 or 1 to indicate input and output cube shapes respectively.

The value of `inp_data_format` and `out_data_format` must be equal.

Note, the fixed-point 8-bit max pool kernel, `xa_nn_maxpool_8` can be used for the quantized int8 datatype.

Function variants available are `xa_nn_maxpool_[p]`, where:

- `[p]`: Input and Output precision in bits

Precision

There are four variants available:

Type	Description
8	8-bit input, 8-bit output
16	16-bit input, 16-bit output
f32	float32 input, float32 output
asym8u	asym8u input, asym8u output

Algorithm

$$z_{h,w,d} = \max(in_{(h*y-stride+i),(w*x-stride+j),d})$$

$$h = 0, \dots, \overline{out-height - 1}, \quad w = 0, \dots, \overline{out-width - 1},$$

$$d = 0, \dots, \overline{out-channels - 1}$$

$$i = 0, \dots, K_H - 1, \quad j = 0, \dots, K_W - 1$$

in denotes padded input cube, *z* denotes output.

K_H, K_W denote `kernel_height`, `kernel_width` respectively.

Prototype

```
WORD32 xa_nn_maxpool_getsize
(WORD32 input_channels, WORD32 inp_precision, WORD32 out_precision,
 WORD32 input_height, WORD32 input_width, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 inp_data_format, WORD32 out_data_format);

WORD32 xa_nn_maxpool_8
(WORD8 * p_out, const WORD8 * p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 out_data_format,
 VOID * p_scratch);
WORD32 xa_nn_maxpool_16
(WORD16 * p_out, const WORD16 * p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 out_data_format,
 VOID * p_scratch);
WORD32 xa_nn_maxpool_f32
(FLOAT32 * p_out, const FLOAT32 * p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 out_data_format,
 VOID * p_scratch);
WORD32 xa_nn_maxpool_asym8u
(UWORD8* p_out, const UWORD8* p_inp, WORD32 input_height,
 WORD32 input_width, WORD32 input_channels, WORD32 kernel_height,
 WORD32 kernel_width, WORD32 x_stride, WORD32 y_stride,
 WORD32 x_padding, WORD32 y_padding, WORD32 out_height,
 WORD32 out_width, WORD32 inp_data_format, WORD32 out_data_format,
 VOID *p_scratch);
```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *, const WORD16 *, const UWORD8 *, const FLOAT32 *	p_inp	input_height * input_width * input_channels	Input cube
WORD32	input_height		Input height
WORD32	input_width		Input width
WORD32	input_channels		Input number of channels
WORD32	kernel_height		Pooling window height
WORD32	kernel_width		Pooling window width
WORD32	x_stride		Horizontal stride over input
WORD32	y_stride		Vertical stride over input
WORD32	x_padding		Left padding width on input
WORD32	y_padding		Top padding height on input

WORD32	out_height		Output height
WORD32	out_width		Output width
WORD32	inp_data_format		Input data format 0:SHAPE_CUBE_DWH_T 1:SHAPE_CUBE_WHD_T
WORD32	out_data_format		Input data format 0:SHAPE_CUBE_DWH_T 1:SHAPE_CUBE_WHD_T
Output			
WORD8 *, WORD16 *, UWORD8 *, FLOAT32 *	p_out	out_height * out_width * input_channels	Output cube
Temporary			
VOID *	p_scratch	xa_nn_maxpool_ getsize()	Temporary / scratch memory

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_inp, p_out	Cannot be NULL
	Should not overlap
	Should not overlap
	Memory size \geq size returned by xa_nn_maxpool_getsize()
input_height, input_width	Greater than or equal to 1
input_channels	Greater than or equal to 1
kernel_height	{1, 2, ..., input_height}
kernel_width	{1, 2, ..., input_width}
x_stride, y_stride	Greater than or equal to 1
x_padding, y_padding	Greater than or equal to 0
out_height, out_width	Greater than or equal to 1
inp_data_format	Can be 0: SHAPE_CUBE_DWH_T or 1: SHAPE_CUBE_WHD_T
out_data_format	Must be equal to inp_data_format

3.5 Fully connected Layer

3.5.1 Fully Connected Kernel

Description

These kernels perform the operation of multiplication of weight matrix with input vectors in a fully connected neural network layer i.e. $z = \text{weight} * \text{input} + \text{bias}$. The column dimension of `weight` must match the row dimension of `input`. Bias and resulting output vector `z` have as many numbers of rows as `weight` matrix.

The `bias_shift` and `acc_shift` arguments are provided in kernel API to adjust Q format of bias and output, respectively. Both `bias_shift` and `acc_shift` can be either positive or negative, where positive value denotes a left shift and negative value denotes a right shift.

`bias_shift` is the shift in number of bits applied to the bias to make it in the same Q format as weight X input multiplication – accumulation result. `acc_shift` is the shift in number of bits applied to the accumulator to obtain the output in desired Q format.

Note: `acc_shift` and `bias_shift` are not relevant in the case of floating point and quantized 8-bit kernels.

For conversion from higher precision accumulator to lower precision output, symmetric rounding is used.

The precision of output is the same as precision of input vector.

The arguments `input_zero_bias`, `weight_zero_bias` are provided to convert the quantized 8-bit inputs into their real values and perform Fully Connected kernel operation. The `out_zero_bias`, `out_multiplier` and `out_shift` values are used to quantize real values of output back to 8-bit.

Function variants available (for fixed point) are `xa_nn_fully_connected_[p]x[q]_[r]`, where:

- `[p]`: Weight matrix precision in bits
- `[q]`: Input vector precision in bits
- `[r]`: Output vector precision in bits

Precision

There are six variants available:

Type	Description
<code>16x16_16</code>	16-bit weight matrix, 16-bit input vector, 16-bit output
<code>8x16_16</code>	8-bit weight matrix, 16-bit input vector, 16-bit output
<code>8x8_8</code>	8-bit weight matrix, 8-bit input vector, 8-bit output
<code>f32</code>	float32 weight matrix, float32 input vector, float32 output
<code>asym8uxasym8u_asym8u</code>	asym8u weight matrix, asym8u input vector, asym8u output
<code>sym8sxsym8s_asym8s</code>	sym8s weight matrix, asym8s input vector, asym8s output

Algorithm

$$z_n = 2^{acc-shift} \left(\sum_{m=0}^{W_D-1} weight_{n,m} \cdot input_m + 2^{bias-shift} bias_n \right),$$

$$n = 0, \dots, \overline{out-depth} - 1$$

where W_D represents weight_depth

For floating point and quantized 8-bit routines, acc_shift=0 and bias_shift=0

Thus, $2^{acc-shift} = 2^{bias-shift} = 1$

Prototype

```
WORD32 xa_nn_fully_connected_16x16_16
(WORD16 * p_out,          WORD16 * p_weight,          WORD16 * p_inp,
 WORD16 * p_bias,          WORD32 weight_depth,          WORD32 out_depth,
 WORD32 acc_shift,          WORD32 bias_shift);
WORD32 xa_nn_fully_connected_8x16_16
(WORD16 * p_out,          WORD8 * p_weight,          WORD16 * p_inp,
 WORD16 * p_bias,          WORD32 weight_depth,          WORD32 out_depth,
 WORD32 acc_shift,          WORD32 bias_shift);
WORD32 xa_nn_fully_connected_8x8_8
(WORD8 * p_out,          WORD8 * p_weight,          WORD8 * p_inp,
 WORD8 * p_bias,          WORD32 weight_depth,          WORD32 out_depth,
 WORD32 acc_shift,          WORD32 bias_shift);
WORD32 xa_nn_fully_connected_f32
(FLOAT32 * p_out,          FLOAT32 * p_weight,          FLOAT32 * p_inp,
 FLOAT32 * p_bias,          WORD32 weight_depth,          WORD32 out_depth);
WORD32 xa_nn_fully_connected_asym8uxasym8u_asym8u
(UWORD8 * p_out,          const UWORD8 * p_weight, const UWORD8 * p_inp,
 const WORD32 * p_bias, WORD32 weight_depth,          WORD32 out_depth,
 WORD32 input_zero_bias, WORD32 weight_zero_bias, WORD32 out_multiplier,
 WORD32 out_shift,          WORD32 out_zero_bias);
WORD32 xa_nn_fully_connected_sym8sxasym8s_asym8s
(WORD8 * p_out,          const WORD8 * p_weight, const WORD8 * p_inp,
 const WORD32 * p_bias, WORD32 weight_depth,          WORD32 out_depth,
 WORD32 input_zero_bias, WORD32 out_multiplier, WORD32 out_shift,
 WORD32 out_zero_bias);
```

Arguments

Type	Name	Size	Description
Input			
WORD16 *, WORD8 *, pFLOAT32, const UWORD8 *, const WORD8 *	p_weight	out_depth* weight_depth	Weight matrix, fixed, floating point, asym8u or sym8s
WORD16 *, WORD8 *, pFLOAT32, const UWORD8 *	p_inp	weight_depth* 1	Input vector, fixed, floating point, asym8u or asym8s

const WORD8 *			
WORD16 *, WORD8 *, pFLOAT32, WORD32 *	p_bias	out_depth*1	Bias vector, fixed or floating point, 32-bit for quantized kernels
WORD32	out_depth		Number of rows in weight matrix, bias and output vector
WORD32	weight_depth		Number of columns in weight matrix and rows in input vector
WORD32	acc_shift		Shift applied to accumulator
WORD32	bias_shift		Shift applied to bias
WORD32	input_zero_bias		Zero offset of input
WORD32	weight_zero_bias		Zero offset of weights
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
Output			
WORD8 *, WORD16 *, pFLOAT32, WORD8 *, UWORD8*	p_out	out_depth*1	Output vector, fixed, floating point, asym8u or asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
weight_depth	Multiple of 4 for fixed point kernels.(2 in case of floating point). No restriction for quantized 8-bit kernels. row_stride >= cols
p_weight, p_inp, p_bias, p_out	Aligned on 16-byte boundary, should not overlap.. (size of one element)-byte boundary in case of floating point and quantized 8-bit kernels).
p_weight, p_inp, p_out	Cannot be NULL
out_depth	Greater than or equal to 1
acc_shift, bias_shift, out_shift	{-31, ...,31}
input_zero_bias	{-255, ...,0} for asym8u, {-127, ...,128} for asym8s
weight_zero_bias	{-255, ...,0} for asym8u
out_zero_bias	{0, ...,255} for asym8u, {-128,,127} for asym8s

3.6 Basic Operations and Miscellaneous Kernels

3.6.1 Interpolation Kernel

Description

This kernel performs interpolation between two input vectors h and y using interpolation factor from vector x to get output vector z .

The interpolation kernel accepts 16-bit inputs and 16-bit interpolation factor in Q15 format and produces 16-bit output in Q15 format.

Precision

Type	Description
16-bit	16-bit input, 16-bit interpolation factor, 16-bit output

Algorithm

$$z_n = x_n * y_n + (1 - x_n) * h_n, \quad n = 0 \dots, \overline{num-elements} - 1$$

x_n represents interpolation factor.

y_n represents first input, h_n represents second input.

z_n represents output.

Prototype

```
WORD32 xa_nn_vec_interpolation_q15
(WORD16 * p_out,          WORD16 * p_ifact,          WORD16 * p_inp1, WORD16 * p_inp2,          WORD32
 num_elements);
```

Arguments

Type	Name	Size	Description
Input			
WORD16 *	p_ifact	num_elements	Interpolation factor vector
WORD16 *	p_inp1	num_elements	First input vector
WORD16 *	p_inp2	num_elements	Second input vector
WORD32	num_elements		Number of elements
Output			
WORD16 *	p_out	num_elements	Output vector

Returns

- 0: no error

- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_ifact, p_inp1, p_inp2, p_out	Aligned on 8-byte boundary
	Should not overlap
	Cannot be NULL
num_elements	Multiple of 4

3.6.2 Dot Product Kernels

Description

These kernels perform the dot product operations between two sets of input vectors `p_inp1` and `p_inp2` to get output vector `p_out`. The supported precisions are: `f32xf32_f32` and `16x16_asym8s`.

Function variants available are `xa_nn_elm_quantize_[p]x[q]_[r]`, where:

- `[p], [q]`: Input precision
- `[r]`: Output precision

Precision

There are two variants available:

Type	Description
<code>f32xf32_f32</code>	float32 input, float32 output
<code>16x16_asym8s</code>	16-bit input, asym8s output

Algorithm

$$z_n = \left(\sum_{m=0}^{vec_length-1} inp1_m \cdot inp2_m + bias_n \right)$$

$$n = 0, \dots, \overline{vec_count} - 1$$

Prototype

```
WORD32 xa_nn_dot_prod_f32xf32_f32
(FLOAT32 * p_out,          const FLOAT32 * p_inp1,          const FLOAT32 * p_inp2,
 WORD32 vec_length,        WORD32 num_vecs);
WORD32 xa_nn_dot_prod_16x16_asym8s
(WORD8 * p_out,            const WORD16 * p_inp1,            const WORD16 * p_inp2,
 const WORD32 * bias_ptr,  WORD32 vec_length,              WORD32 out_multiplier,
 WORD32 out_shift,        WORD32 out_zero_bias,            WORD32 vec_count);
```

Arguments

Type	Name	Size	Description
Input			
const FLOAT32 * const WORD16 *	p_inp1	vec_length	First input vector
const FLOAT32 * const WORD16 *	p_inp2	vec_length	Second input vector
const WORD32 *	bias_ptr	vec_count	Bias vector
WORD32	vec_length		Length of each vector
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
WORD32	num_vecs, vec_count		Number of input vectors
Output			
FLOAT32 * WORD8 *	p_out	num_vecs	Output vector

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions:

Arguments	Restrictions
p_inp1, p_inp2, p_out	Aligned on (size of one element)-byte boundary Cannot be NULL
vec_length, num_vecs	Greater than 0
out_shift	{-31, ..., 31}
out_multiplier	Greater than 0
out_zero_bias	{-128, ..., 127} for out type asym8s

3.6.3 Elementwise Quantize Kernel

Description

These kernels perform the quantization operation of the `p_inp1` input vector elements to get the output vector `p_out`. The kernel is developed in reference to the Quantize operator implementation in Tensorflow Lite Micro.

Function variants available are `xa_nn_elm_quantize_[p]_[q]`, where:

- [p]: Input precision
- [q]: Output precision

Precision

Type	Description
asym16s_asym8s	asym16s input, asym8s output

Prototype

```
WORD32 xa_nn_elm_quantize_asym16s_asym8s
(WORD8 * p_out,          const WORD16 * p_inp,          WORD32 inp_zero_bias,
 WORD32 out_zero_bias,   WORD32 out_shift,              WORD32 out_multiplier,
 WORD32 num_elm);
```

Arguments

Type	Name	Size	Description
Input			
const WORD16 *	p_inp	num_elm	Input vector
WORD32	inp_zero_bias		Zero offset of input
WORD32	out_zero_bias		Zero offset of output
WORD32	out_shift		Shift value of output
WORD32	out_multiplier		Multiplier value of output
WORD32	num_elm		Number of input elements
Output			
WORD8 *	p_out	num_vecs	Output vector

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions:

Arguments	Restrictions
p_inp, p_out	Aligned on (size of one element)-byte boundary
	Cannot be NULL
num_elm	Greater than 0
out_shift	{-31, ..., 31}
out_multiplier	Greater than 0
inp_zero_bias	{-32768.....,32767} for out type asym8s
out_zero_bias	{-128.....,127} for out type asym8s

3.6.4 Basic Kernels

Description

These kernels perform basic elementwise operations on one or two input vectors x and y to get output vector z . The supported operations are: add, subtract, multiply, floor, minimum and maximum. The supported precisions are: 8-bit, float32 and asym8s.

The 8-bit elementwise minimum and maximum kernels can be also used for asym8s datatype.

Function variants available are `xa_nn_[o]_[p]_[q]`, where:

- `[o]`: Operations: `elm_add`, `elm_sub`, `elm_mul`, `elm_floor`, `elm_min`, `elm_max`
- `[p]`: Input Precision in bits- `input1input2` or `input1`
- `[q]`: Output Precision in bits

Precision

Type	Description
<code>f32xf32_f32</code>	float32 input, float32 output
<code>8x8_8</code>	8-bit input, 8-bit output
<code>asym8sxasym8s_asym8s</code>	asym8s input, asym8s output

Algorithm

$$\begin{array}{lll}
 \text{elm_add:} & z_n = x_n + y_n, & n = 0 \dots, \overline{\text{num-elm} - 1} \\
 \text{elm_sub:} & z_n = x_n - y_n, & n = 0 \dots, \overline{\text{num-elm} - 1} \\
 \text{elm_mul:} & z_n = x_n * y_n, & n = 0 \dots, \overline{\text{num-elm} - 1} \\
 \text{elm_floor:} & z_n = \lfloor x_n \rfloor, & n = 0 \dots, \overline{\text{num-elm} - 1} \\
 \text{elm_min:} & z_n = \min(x_n, y_n), & n = 0 \dots, \overline{\text{num-elm} - 1} \\
 \text{elm_max:} & z_n = \max(x_n, y_n), & n = 0 \dots, \overline{\text{num-elm} - 1}
 \end{array}$$

x_n represents first input, y_n represents second input.

z_n represents output.

Prototype

```

WORD32 xa_nn_elm_floor_f32_f32
(FLOAT32 * p_out, const FLOAT32 * p_inp, WORD32 num_elm);
WORD32 xa_nn_elm_add_asym8sxasym8s_asym8s
(WORD8 * p_out, WORD32 out_zero_bias, WORD32 out_shift,
WORD32 out_multiplier, WORD32 out_activation_min, WORD32 out_activation_max,
const WORD8 * p_inpl, WORD32 inpl_zero_bias, WORD32 inpl_shift,
WORD32 inpl_multiplier, const WORD8 * p_inp2, WORD32 inp2_zero_bias,
WORD32 inp2_shift, WORD32 inp2_multiplier, WORD32 left_shift,
WORD32 num_elm);
WORD32 xa_nn_elm_sub_asym8sxasym8s_asym8s
(WORD8 * p_out, WORD32 out_zero_bias, WORD32 out_left_shift,
```

```

WORD32 out_multiplier,      WORD32 out_activation_min,  WORD32 out_activation_max,
const WORD8 * p_inpl,      WORD32 inpl_zero_bias,      WORD32 inpl_left_shift,
WORD32 inpl_multiplier,    const WORD8 * p_inp2,      WORD32 inp2_zero_bias,
WORD32 inp2_left_shift,    WORD32 inp2_multiplier,    WORD32 left_shift,
WORD32 num_elm);
WORD32 xa_nn_elm_mul_asym8sxasym8s_asym8s
(WORD8 * p_out,              WORD32 out_zero_bias,              WORD32 out_shift,
WORD32 out_multiplier,      WORD32 out_activation_min,  WORD32 out_activation_max,
const WORD8 * p_inpl,      WORD32 inpl_zero_bias,      const WORD8 * p_inp2,
WORD32 inp2_zero_bias,      WORD32 num_elm);
WORD32 xa_nn_elm_min_8x8_8
(WORD8* p_out,              const WORD8* p_in1,              const WORD8* p_in2,
WORD32 num_element);
WORD32 xa_nn_elm_max_8x8_8
(WORD8* p_out,              const WORD8* p_in1,              const WORD8* p_in2,
WORD32 num_element);

```

Arguments

Type	Name	Size	Description
Input			
const WORD8 * FLOAT32 *	p_inpl, p_inp, p_inl	num_elm	First input vector
const WORD8 *	p_inp2, P_in2	num_elm	Second input vector
WORD32	num_elm/num_element		Number of elements
WORD32	out_zero_bias		Zero bias of output
WORD32	out_shift		Shift value of output
WORD32	out_multiplier		Multiplier value of output
WORD32	out_activation_min		Activation min of output
WORD32	out_activation_max		Activation max of output
WORD32	inpl_zero_bias		Zero bias of input 1
WORD32	inpl_shift		Shift value of input 1
WORD32	inpl_multiplier		Multiplier value of input 1
WORD32	inp2_zero_bias		Zero bias of input 2
WORD32	inp2_shift		Shift value of input 2
WORD32	inp2_multiplier		Multiplier value of input 2
WORD32	left_shift		Global left shift value for inputs.
Output			
WORD8 * FLOAT32 *	p_out	num_elm	Output vector

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions:

Arguments	Restrictions
p_inpl, p_inp2, p_inp, p_inl, p_in2 p_out	Aligned on (size of one element)-byte boundary Cannot be NULL

num_elm, num_element	Greater than 0
inp1_zero_bias, inp2_zero_bias	{-127....., 128} for asym8s input
inp1_shift, inp2_shift, out_shift	{-31 31} for fixed point and quantized 8-bit APIs
left_shift	{0 31}
inp1_multiplier, inp2_multiplier out_multiplier	Shouldn't be less than 0.
out_zero_bias	{-128....., 127} for asym8s output
out_activation_min, out_activation_max	{-128....., 127} for asym8s output out_activation_min < out_activation_max

3.6.5 Basic Kernels with Broadcasting

Description

These kernels perform a broadcast operation and apply an arithmetic operator. The supported operators are: elementwise minimum and maximum.

Details of the broadcast operation can be found at [Tensorflow Broadcasting semantics](#) [4].

There are two variants of these kernels, one for 4-dimensional and another for 8-dimensional input/output tensors. Input tensors smaller than these dimensions must have their shapes extended^{4.1} to match either of these two.

Tensors must also be broadcast compatible (as these kernels do not perform any runtime checks and depend on the TensorFlow infrastructure)

The input to these kernels are the IO pointers to tensors stored in row-major format, the shape of the resulting broadcasted output and the input 'strides' [5].

Function variants available are `xa_nn_[op]_[d]_Bcast_[p]`, where:

- `[op]`: Operation: `elm_min`, `elm_max`
- `[d]`: Number of IO dimensions: 4D, 8D
- `[p]`: Input/Output precision in bits as `[in1_precision]x[in2_precision]_[out_precision]`

Precision

Type	Description
8x8_8	Signed 8-bit inputs, signed 8-bit output

Algorithm

$$p_out[i_0][i_1] \dots [i_N] = [op](p_in1([i_0 \ i_1 \dots i_N] \cdot [s1_0 \ s1_1 \dots s1_N]), p_in2([i_0 \ i_1 \dots i_N] \cdot [s2_0 \ s2_1 \dots s2_N]))$$

Where,

- $i_n \in (0 \text{ out_extents}[n])$, and, $n \in (0 \ 4]$ for 4D tensors, or, $(0 \ 8]$ for 8D Tensors
- $s1_n = \text{in1_strides}[n]$, with n defined the same as above
- $s2_n = \text{in2_strides}[n]$, with n defined the same as above

Prototypes

```
WORD32 xa_nn_elm_min_4D_Bcast_8x8_8(
    WORD8* __restrict__ p_out,  const int* const out_extents,
    const WORD8* __restrict__ p_in1,  const int* const in1_strides,
    const WORD8* __restrict__ p_in2,  const int* const in2_strides )
WORD32 xa_nn_elm_max_4D_Bcast_8x8_8(
    WORD8* __restrict__ p_out,  const int* const out_extents,
```



```

    const WORD8* __restrict__ p_in1, const int* const in1_strides,
    const WORD8* __restrict__ p_in2, const int* const in2_strides )
WORD32 xa_nn_elm_min_8D_Bcast_8x8_8(
    WORD8* __restrict__ p_out, const int* const out_extents,
    const WORD8* __restrict__ p_in1, const int* const in1_strides,
    const WORD8* __restrict__ p_in2, const int* const in2_strides )
WORD32 xa_nn_elm_max_8D_Bcast_8x8_8(
    WORD8* __restrict__ p_out, const int* const out_extents,
    const WORD8* __restrict__ p_in1, const int* const in1_strides,
    const WORD8* __restrict__ p_in2, const int* const in2_strides )

```

Arguments

Type	Name	Size	Description
Input			
const WORD8*	p_in1	-	First input tensor in row-major
const int* const	in1_strides	4 or 8	Strides for first input tensor
const WORD8*	p_in2	-	Second input tensor in row-major
const int* const	in2_strides	4 or 8	Strides for second input tensor
const int* const	out_extents	4 or 8	Broadcasted output shape
Output			
WORD8*	p_out	prod(out_extents)	Output tensor in row-major

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_in1, p_in2 p_out	Aligned on byte boundary Cannot be NULL
out_extents, in1_strides, in2_strides	Positive integers

3.6.6 Elementwise Comparison Kernels

Description

These kernels perform elementwise comparison operations on two input vectors x and y to get the output vector z . The supported operations are: equal, not equal, greater, greater equal, less, less equal. The output for all the comparison kernels is a boolean value that requires 1-byte space. The supported precisions are: asym8s.

Function variants available are `xa_nn_[o]_[p]`, where:

- [o]: Operations: elm_equal, elm_notequal, elm_greater, elm_greaterequal, elm_less, elm_lessequal
- [p]: Input Precision in bits- input1xinput2

Precision

Type	Description
asym8sxasym8s	asym8s inputs, boolean(1-byte) output

Algorithm

elm_equal: $z_n = (x_n == y_n), \quad n = 0 \dots, \overline{num-elm - 1}$
 elm_notequal: $z_n = (x_n \neq y_n), \quad n = 0 \dots, \overline{num-elm - 1}$
 elm_greater: $z_n = (x_n > y_n), \quad n = 0 \dots, \overline{num-elm - 1}$
 elm_greaterequal: $z_n = (x_n \geq y_n), \quad n = 0 \dots, \overline{num-elm - 1}$
 elm_less: $z_n = (x_n < y_n), \quad n = 0 \dots, \overline{num-elm - 1}$
 elm_lessequal: $z_n = (x_n \leq y_n), \quad n = 0 \dots, \overline{num-elm - 1}$

x_n represents first input, y_n represents second input.

z_n represents output.

Prototype

```

WORD32 xa_nn_elm_equal_asym8sxasym8s
(WORD8 * p_out,          const WORD8 * p_inp1,      WORD32 inp1_zero_bias,
 WORD32 inp1_shift,      WORD32 inp1_multiplier,    const WORD8 * p_inp2,
 WORD32 inp2_zero_bias,  WORD32 inp2_shift,          WORD32 inp2_multiplier,
 WORD32 left_shift,      WORD32 num_elm);
WORD32 xa_nn_elm_notequal_asym8sxasym8s
(WORD8 * p_out,          const WORD8 * p_inp1,      WORD32 inp1_zero_bias,
 WORD32 inp1_shift,      WORD32 inp1_multiplier,    const WORD8 * p_inp2,
 WORD32 inp2_zero_bias,  WORD32 inp2_shift,          WORD32 inp2_multiplier,
 WORD32 left_shift,      WORD32 num_elm);
WORD32 xa_nn_elm_greater_asym8sxasym8s
(WORD8 * p_out,          const WORD8 * p_inp1,      WORD32 inp1_zero_bias,
 WORD32 inp1_shift,      WORD32 inp1_multiplier,    const WORD8 * p_inp2,
 WORD32 inp2_zero_bias,  WORD32 inp2_shift,          WORD32 inp2_multiplier,
 WORD32 left_shift,      WORD32 num_elm);
WORD32 xa_nn_elm_greaterequal_asym8sxasym8s
(WORD8 * p_out,          const WORD8 * p_inp1,      WORD32 inp1_zero_bias,
 WORD32 inp1_shift,      WORD32 inp1_multiplier,    const WORD8 * p_inp2,
 WORD32 inp2_zero_bias,  WORD32 inp2_shift,          WORD32 inp2_multiplier,
 WORD32 left_shift,      WORD32 num_elm);
WORD32 xa_nn_elm_less_asym8sxasym8s
(WORD8 * p_out,          const WORD8 * p_inp1,      WORD32 inp1_zero_bias,
 WORD32 inp1_shift,      WORD32 inp1_multiplier,    const WORD8 * p_inp2,
 WORD32 inp2_zero_bias,  WORD32 inp2_shift,          WORD32 inp2_multiplier,
 WORD32 left_shift,      WORD32 num_elm);
WORD32 xa_nn_elm_lessequal_asym8sxasym8s
(WORD8 * p_out,          const WORD8 * p_inp1,      WORD32 inp1_zero_bias,
 WORD32 inp1_shift,      WORD32 inp1_multiplier,    const WORD8 * p_inp2,
 WORD32 inp2_zero_bias,  WORD32 inp2_shift,          WORD32 inp2_multiplier,
 WORD32 left_shift,      WORD32 num_elm);

```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *	p_inp1	num_elm	First input vector
const WORD8 *	p_inp2	num_elm	Second input vector
WORD32	num_elm		Number of elements
WORD32	inp1_zero_bias		Zero bias of input 1
WORD32	inp1_shift		Shift value of input 1
WORD32	inp1_multiplier		Multiplier value of input 1
WORD32	inp2_zero_bias		Zero bias of input 2
WORD32	inp2_shift		Shift value of input 2
WORD32	inp2_multiplier		Multiplier value of input 2
WORD32	left_shift		Global left shift value for inputs.
Output			
WORD8 *	p_out	num_elm	Output vector

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions:

Arguments	Restrictions
p_inp1, p_inp2, p_out,	Aligned on (size of one element)-byte boundary
	Cannot be NULL
num_elm	Greater than 0
inp1_zero_bias, inp2_zero_bias	{-127, 128} for asym8s input
inp1_shift, inp2_shift	{-31 31} for fixed point and quantized 8-bit APIs
inp1_multiplier, inp2_multiplier	Shouldn't be less than 0.
left_shift	{0 31}

3.6.7 Elementwise Logical Kernels

Description

These kernels perform elementwise logical operations on two boolean input vectors x and y to get the boolean output vector z . The supported operations are: `logical_and`, `logical_or`, `logical_not`. The inputs and output for all the logical kernels are boolean values that requires 1-byte space each. The supported precisions are: `bool`.

Function variants available are `xa_nn_[o]_[p]`, where:

- [o]: Operations: elm_logicaland, elm_logicalor, elm_logicalnot
- [p]: Input Precision in bits- input1xinput2

Precision

Type	Description
boolxbool	boolean(1-byte) inputs, boolean(1-byte) output

Algorithm

elm_logicaland: $z_n = (x_n \&\& y_n), \quad n = 0 \dots, \overline{num_elm - 1}$
 elm_logicalor: $z_n = (x_n \|\| y_n), \quad n = 0 \dots, \overline{num_elm - 1}$
 elm_logicalnot: $z_n = (!x_n), \quad n = 0 \dots, \overline{num_elm - 1}$

x_n represents first input, y_n represents second input.

z_n represents output.

Prototype

```

WORD32 xa_nn_elm_logicaland_boolxbool_bool
(WORD8 * __restrict__ p_out, const WORD8 * __restrict__ p_inp1,
 const WORD8 * __restrict__ p_inp2, WORD32 num_elm);

WORD32 xa_nn_elm_logicalor_boolxbool_bool
(WORD8 * __restrict__ p_out, const WORD8 * __restrict__ p_inp1,
 const WORD8 * __restrict__ p_inp2, WORD32 num_elm);

WORD32 xa_nn_elm_logicalnot_bool_bool
(WORD8 * __restrict__ p_out, const WORD8 * __restrict__ p_inp,
 WORD32 num_elm);

```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *	p_inp1 / p_inp	num_elm	First input vector
const WORD8 *	p_inp2	num_elm	Second input vector
WORD32	num_elm		Number of elements
Output			
WORD8 *	p_out	num_elm	Output vector

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions:

Arguments	Restrictions
p_inp1/p_inp, p_inp2, p_out	Aligned on (size of one element)-byte boundary
	Cannot be NULL
num_elm	Greater than 0

3.6.8 L2 Normalization Kernel

Description

This kernel performs L2 normalization of an input vector x to get output vector z , which means every element of input vector x is divided by L2 norm of x , this gives an output vector z whose L2 norm is 1.

The L2 Normalization kernel accepts asym8s input vector and produces asym8s output vector.

Precision

Type	Description
asym8s	asym8s input, asym8s output

Algorithm

$$z_n = \frac{x_n}{\sqrt{\sum_{n=1}^N |x_n|^2}}, \quad n = 1 \dots, \overline{\text{num-elements}}$$

x_n represents input vector.

z_n represents output vector.

Prototype

```
WORD32 xa_nn_l2_norm_asym8s_asym8s
(WORD8 * p_out, const WORD8 * p_inp, WORD32 zero_point, WORD32 num_elm);
```

Arguments

Type	Name	Size	Description
Input			
const WORD8 *	p_inp	num_elm	Input vector
WORD32	zero_point		Input zero bias
WORD32	num_elm		Number of elements
Output			
WORD8 *	p_out	num_elm	Output vector

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions

Arguments	Restrictions
p_inp, p_out	Aligned on (size of one element)-byte boundary
	Should not overlap
	Cannot be NULL
zero_point	{-128....., 127} for asym8s input
num_elm	Greater than 0

3.6.9 Reduce Kernels

Description

These kernels perform reduction operations on an input vector x based on the dimensions given in axis vector and get the output vector z . The supported operations are: `reduce_max` and `reduce_mean`. The supported precisions are: `asym8s`. The kernels presently support upto 4 dimensions and the input data is assumed to be in “NHWC” or “DWHN” data format (Depth or channels dimension is written first).

Note: The axis vector should have non-duplicate values to avoid larger execution time and poor performance.

For the `reduce_max` kernel, the input and output quantization are expected to be same. Thus, the API does not include quantization specific multiplier, shift and zero bias arguments. For the dimensions mentioned in the axis vector, max operation is carried out thereby reducing the dimension size to 1.

For the `reduce_mean` kernel, the input and output quantization can be different. The arguments `inp_zero_bias`, `out_zero_bias`, `out_multiplier` and `out_shift` are provided for the Mean operation and requantization into `asym8s` output. For the dimensions mentioned in the axis vector, mean operation is carried out thereby reducing the dimension size to 1.

Note: The total number of elements in axis dimensions i.e. the values which are to be reduced should not be more than 127 for the `reduce_mean` kernel.

These kernels require temporary buffer for reduce operation. This temporary buffer is provided by `p_scratch` argument of kernel API. The size of temporary buffer should be queried using `xa_nn_reduce_getsize_nhwc()` helper API. The `reduce_ops` argument accepts an enumerator that will state the reduce operation type. It can take the following values: `REDUCE_MAX` and `REDUCE_MEAN`.

Function variants available are `xa_nn_reduce_[o]_[n]_[p]`, where:

- [o]: Operations: `reduce_max`, `reduce_mean`
- [n]: Number of dimentions: 4D

- [p]: Input Precision in bits- input_output

Precision

Type	Description
asym8s_asym8s	asym8s input, asym8s output

Algorithm

Reduce Max:

- For every dimension \mathbf{r} in axis:

$$Z_{N,H,W,C} = \max(in_{n,h,w,c}[\mathbf{r}_i], in_{n,h,w,c}[\mathbf{r}_j])$$

Where,

- The values of output dimensions(N, H, W, C) if reduced will be equal to 1
- $\mathbf{r} \in$ dimensions along which reduce max is to be performed .
- \mathbf{r}_i and \mathbf{r}_j are the elements in the input shape along the \mathbf{r} dimension.

Reduce Mean:

- For every dimension \mathbf{r} in axis:

$$S_{N,H,W,C} = \text{sum}(in_{n,h,w,c}[\mathbf{r}_i], in_{n,h,w,c}[\mathbf{r}_j])$$
- Then, we compute the mean

$$Z_{N,H,W,C} = \frac{1}{\prod nElem_{\mathbf{r}}} S_{N,H,W,C}$$

Where,

- The values of output dimensions(N, H, W, C) if reduced will be equal to 1
- $\mathbf{r} \in$ dimensions along which reduce mean is to be performed .
- \mathbf{r}_i and \mathbf{r}_j are the elements in the input shape along the \mathbf{r} dimension.
- $\prod nElem_{\mathbf{r}}$ is the product of number of elements in every \mathbf{r} dimension.

$S_{N,H,W,C}$ represents the intermediate reduce sum output required for reduce mean.

$Z_{N,H,W,C}$ represents the reduce operation output and $in_{n,h,w,c}$ represents the input vector.

Prototype

```
WORD32 xa_nn_reduce_getsize_nhwc
(WORD32 inp_precision, const WORD32 *const p_inp_shape, WORD32 num_inp_dims,
 const WORD32 *p_axis, WORD32 num_axis_dims, WORD32 reduce_ops);
```

```
WORD32 xa_nn_reduce_max_4D_asym8s_asym8s
(WORD8 * p_out, const WORD32 *const p_out_shape, const WORD8 * p_inp,
 const WORD32 *const p_inp_shape, const WORD32 * p_axis,
 WORD32 num_out_dims, WORD32 num_inp_dims, WORD32 num_axis_dims,
 pVOID p_scratch_in);
```

```
WORD32 xa_nn_reduce_mean_4D_asym8s_asym8s
```

```
(WORD8 * p_out,          const WORD32 *const p_out_shape, const WORD8 * p_inp,
const WORD32 *const p_inp_shape,      const WORD32 * p_axis,
WORD32 num_out_dims,   WORD32 num_inp_dims,      WORD32 num_axis_dims,
WORD32 inp_zero_bias,  WORD32 out_multiplier,    WORD32 out_shift,
WORD32 out_zero_bias,  pVOID p_scratch_in);
```

Arguments

Type	Name	Size	Description
Input			
const WORD32 *const	p_out_shape	num_out_dims	Output shape vector containing size in each output dimension.
const WORD8 *	p_inp	Product of all dims in p_inp_shape	Input vector, asym8s
const WORD32 *const	p_inp_shape	num_inp_dims	Input shape vector containing size in each input dimension.
const WORD32 *	p_axis	num_axis_dims	Axis vector, contains dimensions for reduce operation
WORD32	num_out_dims		Number of output dimension
WORD32	num_inp_dims		Number of input dimension
WORD32	num_axis_dims		Number of axis dimension
WORD32	inp_zero_bias		Zero offset of input
WORD32	out_multiplier		Multiplier value of output
WORD32	out_shift		Shift value of output
WORD32	out_zero_bias		Zero offset of output
pVOID	p_scratch	xa_nn_reduce_getsize_nhwc()	Scratch memory pointer
Output			
WORD8 *	p_out	Product of all dims in p_out_shape	Output vector, asym8s

Returns

- 0: no error
- -1: error, invalid parameters

Restrictions:

Arguments	Restrictions
reduce_ops	Should be REDUCE_MAX or REDUCE_MEAN.
p_inp, p_axis, p_out, p_inp_shape, p_out_shape	Aligned on (size of one element)-byte boundary
	Cannot be NULL and cannot overlap
num_inp_dims, num_out_dims, num_axis_dims	Should be more than 0 and less than equal to 4.
p_axis	Shouldn't be less than 0 and more than 4.
	The axis values should be between 0 and (num_inp_dims - 1).
p_inp_shape, p_out_shape	The shape values should be greater than 0.

p_out_shape	The output length i.e the product of all the shape values must be <= 127.
inp_zero_bias out_zero_bias	{-128.....,127} for asym8s
out_multiplier	Greater than 0
out_shift	{-31, ..., 31}

4. HiFi 5 NN Library – Layers

This section explains the APIs of each layer implementation in the NN library. All the layers conform to the “generic NN Layer API” and flow explained in Section 2.

The NN library is a single archive containing all layers and low-level kernels implementations. Each layer has its own header file that defines the APIs specific to the layer. The following sections explain each layer in detail.

Note This version of the library supports GRU, LSTM, and CNN layers.

4.1 GRU Layer

The GRU APIs are defined in `xa_nnlb_gru_api.h`. Refer to the overall signal flow diagram of GRU in [\[1\]](#).

4.1.1 GRU Layer Specification

GRU layer implements the following input-output equations [\[1\]](#):

$$\begin{aligned} z_t &= \text{sigmoid}(W_z * x_t + U_z * \text{prev-h} + b_z) \\ r_t &= \text{sigmoid}(W_r * x_t + U_r * \text{prev-h} + b_r) \\ g &= \tanh(W_h * x_t + U_h * (r_t \cdot \text{prev-h}) + b_h) \\ y_t = h_t &= z_t \cdot g + (1 - z_t) \cdot \text{prev-h} \\ \text{prev-h} &= h_t \end{aligned}$$

x_t : input vector

y_t, h_t : output vector

W, U : weight matrices

prev-h : previous output vector

z_t : update gate vector

r_t : reset gate vector

b : bias vectors

4.1.2 Error Codes Specific to GRU

Other than common error codes explained in Section 2.3, the GRU layer may also report the following error codes, which may be generated during the initialization stage.

- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_IN_FEATS3`
Number of input features is not supported
- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_OUT_FEATS`
Number of output features is not supported
- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_PRECISION`
I/O precision is not supported
- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_COEFF_QFORMAT`
Number of fractional bits for coefficients is not supported.
- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_IO_QFORMAT`
Number of fractional bits for input-output is not supported.
- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_MEMBANK_PADDING`
Membank padding should be 0 or 1.
- `XA_NNLIB_GRU_CONFIG_FATAL_INVALID_PARAM_ID`
Parameter identifier (param_id) is not valid

The following error codes may be generated during the execution stage.

- `XA_NNLIB_GRU_EXECUTE_FATAL_INSUFFICIENT_DATA`
Input data passed in is insufficient
- `XA_NNLIB_GRU_EXECUTE_FATAL_INSUFFICIENT_OUTPUT_BUFFER_SPACE`
Output Buffer Size is not sufficient

³ FEATS := features

4.1.3 API Functions Specific to GRU

4.1.3.1 Query Functions

Table 4-1 GRU Get Persistent Size Function

Function	<code>xa_nnlb_gru_get_persistent_fast</code>
Syntax	<pre>Int32 xa_nnlb_gru_get_persistent_fast(xa_nnlb_gru_init_config_t *config)</pre>
Description	Returns persistent memory size in bytes required by GRU layer.
Parameters	Input: <code>config</code> Initial configuration parameters (see Table 4-7).
Errors	<p>If return value is less than 0, then it is an error. Following are the possible error codes:</p> <ul style="list-style-type: none"> ■ <code>XA_NNLB_FATAL_MEM_ALLOC</code> ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_IN_FEATS</code> Number of input features is not supported ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_OUT_FEATS</code> Number of output features is not supported ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_COEFF_QFORMAT</code> Number of fractional bits for coefficients is not supported. ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_IO_QFORMAT</code> Number of fractional bits for input-output is not supported.

Table 4-2 GRU Get Scratch Size Function

Function	<code>xa_nnlib_gru_get_scratch_fast</code>
Syntax	<pre>Int32 xa_nnlib_gru_get_scratch_fast(xa_nnlib_gru_init_config_t *config)</pre>
Description	Returns scratch memory size in bytes required by GRU layer.
Parameters	Input: <code>config</code> Initial configuration parameters (see Table 4-7).
Errors	<p>If return value is less than 0, then it is an error. Following are the possible error codes:</p> <ul style="list-style-type: none"> ■ <code>XA_NNLIB_FATAL_MEM_ALLOC</code> ■ <code>XA_NNLIB_GRU_CONFIG_FATAL_INVALID_IN_FEATS</code> Number of input features is not supported ■ <code>XA_NNLIB_GRU_CONFIG_FATAL_INVALID_OUT_FEATS</code> Number of output features is not supported ■ <code>XA_NNLIB_GRU_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported ■ <code>XA_NNLIB_GRU_CONFIG_FATAL_INVALID_COEFF_QFORMAT</code> Number of fractional bits for coefficients is not supported ■ <code>XA_NNLIB_GRU_CONFIG_FATAL_INVALID_IO_QFORMAT</code> Number of fractional bits for input-output is not supported

4.1.3.2 Initialization Stage

Table 4-3 GRU Init Function

Function	<code>xa_nnlb_gru_init</code>
Syntax	<pre>Int32 xa_nnlb_gru_init (xa_nnlb_handle_t handle, xa_nnlb_gru_init_config_t *config)</pre>
Description	Reset the GRU Layer API handle into its initial state. Set up the GRU Layer to the specified initial configuration parameters. This function sets <code>prev_h</code> vector to 0; the user can put the desired values in <code>prev_h</code> by using set config <code>XA_NNLB_GRU_RESTORE_CONTEXT</code> (refer to Table 4-11 for more information).
Parameters	<p>Input: <code>handle</code> Pointer to the component persistent memory. This is the opaque handle. Required size: see <code>xa_nnlb_gru_get_persistent_fast</code>. Required alignment: 8 bytes.</p> <p>Input: <code>config</code> Initial configuration parameters (see Table 4-7). Note that the initial configuration parameters <i>must</i> be identical to those passed to query functions.</p>
Errors	<p>If the return value is not <code>XA_NNLB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA_NNLB_FATAL_MEM_ALLOC</code> One of the pointers is invalid. ■ <code>XA_NNLB_FATAL_MEM_ALIGN</code> One of the pointers is not properly aligned. ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_IN_FEATS</code> Number of input features is not supported ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_OUT_FEATS</code> Number of output features is not supported ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported. ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_COEFF_QFORMAT</code> Number of fractional bits for coefficients is not supported. ■ <code>XA_NNLB_GRU_CONFIG_FATAL_INVALID_IO_QFORMAT</code> Number of fractional bits for input-output is not supported.

4.1.3.3 Execution Stage

Table 4-4 GRU Execution Function

Function	<code>xa_nnl-lib-gru-process</code>
Syntax	<pre> Int32 xa_nnl-lib-gru-process(xa_nnl-lib-handle_t handle, void *scratch, void *input, void *output, xa_nnl-lib-shape_t *p_in_shape, xa_nnl-lib-shape_t *p_out_shape) </pre>
Description	Processes one input shape to generate one output shape.
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>scratch</code> A pointer to the scratch buffer. Required alignment: 8 bytes.</p> <p>Input: <code>input</code> A pointer to the input buffer. Input buffer contains input data. Required alignment: 8 bytes.</p> <p>Output: <code>output</code> A pointer to the output buffer. Output is written to output buffer. Required alignment: 8 bytes.</p> <p>Input/Output: <code>p_in_shape</code> Pointer to the shape containing input buffer dimensions. Contains the length of input data passed to GRU layer. Required alignment: 4 bytes.</p> <p>Input/Output: <code>p_out_shape</code> Pointer to the shape for output buffer dimensions. On return, <code>*p_out_shape</code> is filled with the length of output generated by HiFi GRU Layer. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA_NNL-LIB-NO-ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA_NNL-LIB-FATAL-MEM-ALLOC</code> One of the pointers is NULL. ■ <code>XA_NNL-LIB-FATAL-MEM-ALIGN</code> One of the pointers is not properly aligned.

	<ul style="list-style-type: none"> ■ XA_NNLIB_FATAL_INVALID_SHAPE Either input or output shape is invalid. ■ XA_NNLIB_GRU_EXECUTE_FATAL_INSUFFICIENT_DATA Input data passed in insufficient. ■ XA_NNLIB_GRU_EXECUTE_FATAL_INSUFFICIENT_OUTPUT_BUFFER_SPACE Output buffer size is not sufficient.
--	--

Table 4-5 GRU Set Parameter Function Details

Function	<code>xa_nnl-lib-gru-set-config</code>
Syntax	<pre>Int32 xa_nnl-lib-gru-set-config (xa_nnl-lib-handle_t handle, xa_nnl-lib-gru-param-id_t param_id, void *params)</pre>
Description	Sets the parameter specified by <code>param_id</code> to the value passed in the buffer pointed to by <code>params</code> .
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>param_id</code> Identifies the parameter to be written. Refer to Table 4-11 for the list of supported parameters.</p> <p>Input: <code>params</code> A pointer to a buffer that contains the parameter value. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA_NNLIB_NO_ERROR</code>, it implies that function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ XA_NNLIB_FATAL_MEM_ALLOC One of the pointers (<code>handle</code> or <code>params</code>) is NULL. ■ XA_NNLIB_FATAL_MEM_ALIGN One of the pointers (<code>handle</code> or <code>params</code>) is not aligned correctly. ■ XA_NNLIB_GRU_CONFIG_FATAL_INVALID_PARAM_ID Parameter identifier (<code>param_id</code>) is not valid.

Table 4-6 GRU Get Parameter Function Details

Function	<code>xa_nnlib_gru_get_config</code>
Syntax	<pre>Int32 xa_nnlib_gru_get_config (xa_nnlib_handle_t handle, xa_nnlib_gru_param_id_t param_id, void *params)</pre>
Description	Gets the value of the parameter specified by <code>param_id</code> in the buffer pointed to by <code>params</code> .
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>param_id</code> Identifies the parameter to be read. Refer to Table 4-11 for the list of supported parameters.</p> <p>Output: <code>params</code> A pointer to a buffer that is filled with the parameter value when the function returns. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA>NNLIB_NO_ERROR</code>, it implies that function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA>NNLIB_FATAL_MEM_ALLOC</code> One of the pointers (<code>handle</code> or <code>params</code>) is <code>NULL</code>. ■ <code>XA>NNLIB_FATAL_MEM_ALIGN</code> One of the pointers (<code>handle</code> or <code>params</code>) is not aligned correctly. ■ <code>XA>NNLIB_GRU_CONFIG_FATAL_INVALID_PARAM_ID</code> Parameter identifier (<code>param_id</code>) is not valid.

4.1.4 Structures Specific to GRU

Table 4-7 GRU Config Structure xa_nnlib_gru_init_config_t

Element Type	Element Name	Range	Default	Description
Int32	in_feats	4-2048	256	Number of input features (must be multiple of 4)
Int32	out_feats	4-2048	256	Number of output features (must be multiple of 4)
Int32	pad	0, 1	1	Padding 16 bytes for HiFi 5
Int32	mat_prec	8, 16	16	Matrix input precision
Int32	vec_prec	16	16	Vector input precision
xa_nnlib_gru_precision_t	precision	XA_NNLIB_GRU_16b16b, XA_NNLIB_GRU_8b16b	XA_NNLIB_GRU_16b16b	Coef and I/O precision. Note: Current library supports only 16b16b and 8b16b precision for GRU
Int16	coeff_Qformat	0-15	15	Number of fractional bits for weights and biases
Int16	io_Qformat	0-15	12	Number of fractional bits for input and output

Table 4-8 xa_nnlib_gru_weights_t Parameter Type

Element Type	Element Name	Range	Default	Description
coeff_t *	w_z	NA	NA	Pointer to coefficient matrix w_z.
xa_nnlib_shape_t	shape_w_z	NA	NA	Shape information about w_z.
coeff_t *	u_z	NA	NA	Pointer to coefficient matrix u_z.
xa_nnlib_shape_t	shape_u_z	NA	NA	Shape information about u_z.
coeff_t *	w_r	NA	NA	Pointer to coefficient matrix w_r.
xa_nnlib_shape_t	shape_w_r	NA	NA	Shape information about w_r.
coeff_t *	u_r	NA	NA	Pointer to coefficient matrix u_r.
xa_nnlib_shape_t	shape_u_r	NA	NA	Shape information about u_r.
coeff_t *	w_h	NA	NA	Pointer to coefficient matrix w_h.
xa_nnlib_shape_t	shape_w_h	NA	NA	Shape information about w_h.
coeff_t *	u_h	NA	NA	Pointer to coefficient matrix u_h.
xa_nnlib_shape_t	shape_u_h	NA	NA	Shape information about u_h.

Table 4-9 xa_nnlib_gru_biases_t Parameter Type

Element Type	Element Name	Range	Default	Description
coeff_t *	b_z	NA	NA	Pointer to coefficient matrix b_z.
xa_nnlib_shape_t	shape_b_z	NA	NA	Shape information about b_z.
coeff_t *	b_r	NA	NA	Pointer to coefficient matrix b_r.
xa_nnlib_shape_t	shape_b_r	NA	NA	Shape information about b_r.
coeff_t *	b_h	NA	NA	Pointer to coefficient matrix b_h.
xa_nnlib_shape_t	shape_b_h	NA	NA	Shape information about b_h.

Note GRU requires all weight matrices' and bias vectors' pointers to be 8 bytes aligned.

4.1.5 Enums Specific to GRU

Table 4-10 Enum xa_nnlib_gru_precision_t

Element	Description
XA_NNLIB_GRU_16bx16b	Coef: 16 bits, I/O: 16 bits Fixed Point
XA_NNLIB_GRU_8bx16b	Coef: 8 bits, I/O: 16 bits Fixed Point
XA_NNLIB_GRU_8bx8b	Not supported
XA_NNLIB_flt16xflt16	Not supported

Note Currently, GRU only supports XA_NNLIB_GRU_16bx16b, XA_NNLIB_GRU_8bx16b precision setting.

Table 4-11 describes parameter IDs for parameters supported by GRU. It contains the following columns:

- Parameter ID: Parameter identifier (`param_id`).
- Value type: A pointer (`params`) to a variable of this type is to be passed.
- RW: Indicates whether the parameter can be read (`get`) and/or written (`set`).
- Range: Indicates valid values of the parameter.
- Default: Default value of the parameter
- Description: Brief description of the parameter.

Table 4-11 GRU Specific Parameters

Parameter ID	Value Type	RW	Range	Default	Description
XA_NNLIB_GRU_RESTORE_CONTEXT	vect_t []	RW	NA	NA	Set previous output. This can be used to set prev_h to specific context (size should be equal to number of output features). Upon set config, the buffer passed is copied to persistent memory; upon get config, it returns the prev_h state in the given buffer.
XA_NNLIB_GRU_WEIGHT	xa_nnlib_gru_weights_t	RW	NA	NA	Weight matrices, pointers to weight matrices along with shape information must be passed via xa_nnlib_gru_weights_t structure for set config. Upon get config, it returns pointers to weight matrices along with their shape information in same structure.
XA_NNLIB_GRU_BIAS	xa_nnlib_gru_biases_t	RW	NA	NA	Bias vectors, pointers to bias vectors along with shape information must be passed via xa_nnlib_gru_biases_t structure for set config. Upon get config, it returns pointers to bias vectors along with their shape information in same structure.
XA_NNLIB_GRU_INPUT_SHAPE	xa_nnlib_shape_t	R	NA	NA	Input shape information, get information of the input shape expected by the layer.
XA_NNLIB_GRU_OUTPUT_SHAPE	xa_nnlib_shape_t	R	NA	NA	Output shape information, get information of the output shape expected by layer.

4.2 LSTM Layer

The LSTM APIs are defined in `xa_nnlib_lstm_api.h`.

4.2.1 LSTM Layer Specification

The LSTM layer implements the following forward path input-output equations:

$$\begin{aligned} f_f &= \text{sigmoid}(w_{xf} * \text{frame}_f + \text{prev-h} * w_{hf} + b_f) \\ i_f &= \text{sigmoid}(w_{xi} * \text{frame}_f + \text{prev-h} * w_{hi} + b_i) \\ c\text{-hat}_f &= \tanh(w_{xc} * \text{frame}_f + \text{prev-h} * w_{hc} + b_c) \\ c_f &= f_f * \text{prev-c} + i_f * c\text{-hat}_f \\ o_f &= \text{sigmoid}(w_{xo} * \text{frame}_f + \text{prev-h} * w_{ho} + b_o) \\ h_f &= o_f * \tanh(c_f) \end{aligned}$$

i_f : input gate

h_t : output vector

$c\text{-hat}_f$: intermediate cell state vector

f_f : forget gate

frame_f : Input vector

w_x : weight matrices of input connections

prev-h : previous output vector

prev-c : previous cell output

b : bias vectors

o_f : output gate

c_f : cell state vector

w_h : weight matrices of recurrent connections

4.2.2 Error Codes Specific to LSTM

Other than common error codes explained in Section 2.3, the LSTM layer may also report the following error codes, which may be generated during the initialization stage:

- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_IN_FEATS4`
Number of input features is not supported
- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_OUT_FEATS`
Number of output features is not supported
- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_PRECISION`
I/O precision is not supported
- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_COEFF_QFORMAT`
Number of fractional bits for coefficients is not supported.
- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_CELL_QFORMAT`
Number of fractional bits for cells is not supported

⁴ FEATS: = features

- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_IO_QFORMAT`
Number of fractional bits for input-output is not supported.
- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_MEMBANK_PADDING`
Membank padding should be 0 or 1.
- `XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_PARAM_ID`
Parameter identifier (param_id) is not valid

The following error codes may be generated during the execution stage.

- `XA_NNLIB_LSTM_EXECUTE_FATAL_INSUFFICIENT_DATA`
Input data passed in insufficient
- `XA_NNLIB_LSTM_EXECUTE_FATAL_INSUFFICIENT_OUTPUT_BUFFER_SPACE`
Output Buffer Size is not sufficient

4.2.3 API Functions Specific to LSTM

4.2.3.1 Query Functions

Table 4-12 LSTM Get Persistent Size Function

Function	<code>xa_nnlstm_get_persistent_fast</code>
Syntax	<pre>Int32 xa_nnlstm_get_persistent_fast (xa_nnlstm_init_config_t *config)</pre>
Description	Returns persistent memory size in bytes required by LSTM layer.
Parameters	Input: <code>config</code> Initial configuration parameters (see Table 4-18).
Errors	<p>If return value is less than 0 then it is an error. Following are the possible error codes:</p> <ul style="list-style-type: none"> ■ <code>XA_NNLSTM_FATAL_MEM_ALLOC</code> ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_IN_FEATS</code> Number of input features is not supported ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_OUT_FEATS</code> Number of output features is not supported ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_COEFF_QFORMAT</code> Number of fractional bits for coefficients is not supported. ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_CELL_QFORMAT</code> Number of fractional bits for cells is not supported ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_IO_QFORMAT</code> Number of fractional bits for input-output is not supported. ■ <code>XA_NNLSTM_CONFIG_FATAL_INVALID_MEMBANK_PADDING</code> Membank padding should be 0 or 1.

Table 4-13 LSTM Get Scratch Size Function

Function	<code>xa_nnlib_lstm_get_scratch_fast</code>
Syntax	<code>Int32 xa_nnlib_lstm_get_scratch_fast (xa_nnlib_lstm_init_config_t *config)</code>
Description	Returns scratch memory size in bytes required by LSTM layer.
Parameters	Input: <code>config</code> Initial configuration parameters (see Table 4-18).
Errors	<p>If return value is less than 0 then it is an error, the possible error codes are:</p> <ul style="list-style-type: none"> ■ <code>XA_NNLIB_FATAL_MEM_ALLOC</code> ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_IN_FEATS</code> Number of input features is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_OUT_FEATS</code> Number of output features is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_COEFF_QFORMAT</code> Number of fractional bits for coefficients is not supported. ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_CELL_QFORMAT</code> Number of fractional bits for cells is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_IO_QFORMAT</code> Number of fractional bits for input-output is not supported. ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_MEMBANK_PADDING</code> Membank padding should be 0 or 1.

4.2.3.2 Initialization Stage

Table 4-14 LSTM Init Function

Function	<code>xa_nnlib_lstm_init</code>
Syntax	<pre>Int32 xa_nnlib_lstm_init (xa_nnlib_handle_t handle, xa_nnlib_lstm_init_config_t *config)</pre>
Description	<p>Reset the LSTM layer API handle into its initial state. Set up the LSTM layer to the specified initial configuration parameters. This function sets <code>prev_h</code> vector and <code>prev_c</code> vector to 0; the user can put the desired values in <code>prev_h</code> and <code>prev_c</code> by using set config <code>XA_NNLIB_LSTM_RESTORE_CONTEXT_OUTPUT</code> and <code>XA_NNLIB_LSTM_RESTORE_CONTEXT_CELL</code> respectively (refer to Table 4-22 for more information).</p>
Parameters	<p>Input: <code>handle</code> Pointer to the component persistent memory. This is the opaque handle. Required size: see <code>xa_nnlib_lstm_get_persistent_fast</code>. Required alignment: 8 bytes.</p> <p>Input: <code>config</code> Initial configuration parameters (see Table 4-18). Note that the initial configuration parameters MUST be identical to those passed to query functions.</p>
Errors	<p>If the return value is not <code>XA_NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA_NNLIB_FATAL_MEM_ALLOC</code> One of the pointers is invalid. ■ <code>XA_NNLIB_FATAL_MEM_ALIGN</code> One of the pointers is not properly aligned. ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_IN_FEATS</code> Number of input features is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_OUT_FEATS</code> Number of output features is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported ■ <code>XA_NNLIB_LSTM_CONFIG_FATAL_INVALID_COEFF_QFORMAT</code> Number of fractional bits for coefficients is not supported.

- **XA>NNLIB_LSTM_CONFIG_FATAL_INVALID_CELL_QFORMAT**
Number of fractional bits for cells is not supported
- **XA>NNLIB_LSTM_CONFIG_FATAL_INVALID_IO_QFORMAT**
Number of fractional bits for input-output is not supported
- **XA>NNLIB_LSTM_CONFIG_FATAL_INVALID_MEMBANK_PADDING**
Membank padding should be 0 or 1.

4.2.3.3 Execution Stage

Table 4-15 LSTM Execution Function

Function	<code>xa_nnlib_lstm_process</code>
Syntax	<pre> Int32 xa_nnlib_lstm_process (xa_nnlib_handle_t handle, void *scratch, void *input, void *output, xa_nnlib_shape_t *p_in_shape, xa_nnlib_shape_t *p_out_shape) </pre>
Description	Processes one input shape to generate one output shape.
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>scratch</code> A pointer to the scratch buffer. Required alignment: 8 bytes.</p> <p>Input: <code>input</code> A pointer to the input buffer. Input buffer contains input data. Required alignment: 8 bytes.</p> <p>Output: <code>output</code> A pointer to the output buffer. Output is written to the output buffer. Required alignment: 8 bytes.</p> <p>Input/Output: <code>p_in_shape</code> Pointer to the shape containing input buffer dimensions. Contains the length of input data passed to LSTM layer. Required alignment: 4 bytes.</p> <p>Input/Output: <code>p_out_shape</code></p>

	<p>Pointer to the shape for output buffer dimensions. On return, *p_out_shape is filled with the length of output generated by HiFi LSTM layer.</p> <p>Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not XA_NNLIB_NO_ERROR, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none">■ XA_NNLIB_FATAL_MEM_ALLOC One of the pointers is NULL.■ XA_NNLIB_FATAL_MEM_ALIGN One of the pointers is not having proper alignment.■ XA_NNLIB_FATAL_INVALID_SHAPE Either input or output shape is invalid.■ XA_NNLIB_LSTM_EXECUTE_FATAL_INSUFFICIENT_DATA Input data passed in insufficient■ XA_NNLIB_LSTM_EXECUTE_FATAL_INSUFFICIENT_OUTPUT_BUFFER_SPACE Output Buffer Size is not sufficient

Table 4-16 LSTM Set Parameter Function Details

Function	<code>xa_nnlib_lstm_set_config</code>
Syntax	<pre> Int32 xa_nnlib_lstm_set_config (xa_nnlib_handle_t handle, xa_nnlib_lstm_param_id_t param_id, void *params) </pre>
Description	Sets the parameter specified by <code>param_id</code> to the value passed in the buffer pointed to by <code>params</code> .
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>param_id</code> Identifies the parameter to be written. Refer to Table 4-11 for the list of supported parameters.</p> <p>Input: <code>params</code> A pointer to a buffer that contains the parameter value. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA>NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA>NNLIB_FATAL_MEM_ALLOC</code> One of the pointers (<code>handle</code> or <code>params</code>) is <code>NULL</code>. ■ <code>XA>NNLIB_FATAL_MEM_ALIGN</code> One of the pointers (<code>handle</code> or <code>params</code>) is not aligned correctly. ■ <code>XA>NNLIB_LSTM_CONFIG_FATAL_INVALID_PARAM_ID</code> Parameter identifier (<code>param_id</code>) is not valid.

Table 4-17 LSTM Get Parameter Function Details

Function	<code>xa_nnlib_lstm_get_config</code>
Syntax	<pre>Int32 xa_nnlib_lstm_get_config (xa_nnlib_handle_t handle, xa_nnlib_lstm_param_id_t param_id, void *params)</pre>
Description	Gets the value of the parameter specified by <code>param_id</code> in the buffer pointed to by <code>params</code> .
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>param_id</code> Identifies the parameter to be read. Refer to Table 4-11 for the list of supported parameters.</p> <p>Output: <code>params</code> A pointer to a buffer that is filled with the parameter value when the function returns. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA>NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA>NNLIB_FATAL_MEM_ALLOC</code> One of the pointers (<code>handle</code> or <code>params</code>) is <code>NULL</code>. ■ <code>XA>NNLIB_FATAL_MEM_ALIGN</code> One of the pointers (<code>handle</code> or <code>params</code>) is not aligned correctly. ■ <code>XA>NNLIB_LSTM_CONFIG_FATAL_INVALID_PARAM_ID</code> Parameter identifier (<code>param_id</code>) is not valid.

4.2.4 Structures Specific to LSTM

Table 4-18 LSTM Config Structure `xa_nnlib_lstm_init_config_t`

Element Type	Element Name	Range	Default	Description
Int32	<code>in_feats</code>	4-2048	256	Number of input features (must be multiple of 4)
Int32	<code>out_feats</code>	4-2048	256	Number of output features (must be multiple of 4)
Int32	<code>pad</code>	0, 1	1	Padding 16 bytes for HiFi 5 DSP
Int32	<code>mat_prec</code>	8, 16	16	Matrix input precision
Int32	<code>vec_prec</code>	16	16	Vector input precision
<code>xa_nnlib_lstm_precision_t</code>	<code>precision</code>	XA_NNLIB_LSTM_16bx16b, XA_NNLIB_LSTM_8bx16b	XA_NNLIB_LSTM_16bx16b	Coef and I/O precision. Note: The current library supports only 16bx16b and 8bx16b precision for LSTM.
Int16	<code>coeff_Qformat</code>	0-15	15	Number of fractional bits for weights and biases
Int16	<code>cell_Qformat</code>	0-26		Number of fractional bits for cells.
Int16	<code>io_Qformat</code>	0-15	12	Number of fractional bits for input and output

Table 4-19 `xa_nnlib_lstm_weights_t` Parameter Type

Element Type	Element Name	Range	Default	Description
<code>coeff_t *</code>	<code>w_xf</code>	NA	NA	Pointer to coefficient matrix <code>w_xf</code> .
<code>xa_nnlib_shape_t</code>	<code>shape_w_xf</code>	NA	NA	Shape information about <code>w_xf</code> .
<code>coeff_t *</code>	<code>w_xi</code>	NA	NA	Pointer to coefficient matrix <code>w_xi</code> .
<code>xa_nnlib_shape_t</code>	<code>shape_w_xi</code>	NA	NA	Shape information about <code>w_xi</code> .
<code>coeff_t *</code>	<code>w_xc</code>	NA	NA	Pointer to coefficient matrix <code>w_xc</code> .
<code>xa_nnlib_shape_t</code>	<code>shape_w_xc</code>	NA	NA	Shape information about <code>w_xc</code> .
<code>coeff_t *</code>	<code>w_xo</code>	NA	NA	Pointer to coefficient matrix <code>w_xo</code> .
<code>xa_nnlib_shape_t</code>	<code>shape_w_xo</code>	NA	NA	Shape information about <code>w_xo</code> .
<code>coeff_t *</code>	<code>w_hf</code>	NA	NA	Pointer to coefficient matrix <code>w_hf</code> .
<code>xa_nnlib_shape_t</code>	<code>shape_w_hf</code>	NA	NA	Shape information about <code>w_hf</code> .
<code>coeff_t *</code>	<code>w_hi</code>	NA	NA	Pointer to coefficient matrix <code>w_hi</code> .

Element Type	Element Name	Range	Default	Description
xa_nnlb_shape_t	shape_w_hi	NA	NA	Shape information about w_hi.
coeff_t *	w_hc	NA	NA	Pointer to coefficient matrix w_hc.
xa_nnlb_shape_t	shape_w_hc	NA	NA	Shape information about w_hc.
coeff_t *	w_ho	NA	NA	Pointer to coefficient matrix w_ho.
xa_nnlb_shape_t	shape_w_ho	NA	NA	Shape information about w_ho.

Table 4-20 xa_nnlb_lstm_biases_t Parameter Type

Element Type	Element Name	Range	Default	Description
coeff_t *	b_f	NA	NA	Pointer to coefficient matrix b_f.
xa_nnlb_shape_t	shape_b_f	NA	NA	Shape information about b_f.
coeff_t *	b_i	NA	NA	Pointer to coefficient matrix b_i.
xa_nnlb_shape_t	shape_b_i	NA	NA	Shape information about b_i.
coeff_t *	b_c	NA	NA	Pointer to coefficient matrix b_c.
xa_nnlb_shape_t	shape_b_c	NA	NA	Shape information about b_c.
coeff_t *	b_o	NA	NA	Pointer to coefficient matrix b_o.
xa_nnlb_shape_t	shape_b_o	NA	NA	Shape information about b_o.

Note LSTM requires all weight matrices' and bias vectors' pointers to be 8 bytes aligned.

4.2.5 Enums Specific to LSTM

Table 4-21 Enum xa_nnlb_lstm_precision_t

Element	Description
XA_NNLB_LSTM_16bx16b	Coef: 16 bits, I/O: 16 bits Fixed Point
XA_NNLB_LSTM_8bx16b	Coef: 8 bits, I/O: 16 bits Fixed Point
XA_NNLB_LSTM_8bx8b	Not supported
XA_NNLB_flt16xflt16	Not supported

Note Currently, LSTM only supports the XA_NNLB_LSTM_16bx16b, XA_NNLB_LSTM_8bx16b precision setting.

Table 4-22 describes parameter IDs for parameters supported by LSTM. It contains the following columns:

- **Parameter ID:** Parameter identifier (`param_id`).
- **Value type:** A pointer (`params`) to a variable of this type is to be passed.
- **RW:** Indicates whether the parameter can be read (`get`) and/or written (`set`).
- **Range:** Indicates valid values of the parameter.
- **Default:** Default value of the parameter.
- **Description:** Brief description of the parameter.

Table 4-22 LSTM Specific Parameters

Parameter ID	Value Type	RW	Range	Default	Description
XA_NNLIB_LSTM_RESTORE_CONTEXT_OUTPUT	<code>vect_t []</code>	RW	NA	NA	Set previous output. This can be used to set <code>prev_h</code> to specific context (size should be equal to number of output features). Upon set config, the buffer passed is copied to persistent memory; upon get config, it returns the <code>prev_h</code> state in the given buffer.
XA_NNLIB_LSTM_RESTORE_CONTEXT_CELL	<code>vect_t []</code>	RW	NA	NA	Set previous cell state. This can be used to set <code>prev_c</code> to specific cell context (size should be equal to number of output features). Upon set config, the buffer passed is copied to persistent memory; upon get config, it returns the <code>prev_c</code> state in the given buffer.
XA_NNLIB_LSTM_WEIGHT	<code>xa_nnl-lib_lstm_weights_t</code>	RW	NA	NA	Weight matrices, pointers to weight matrices along with shape information needs to be passed via <code>xa_nnl-lib_lstm_weights_t</code> structure for set config. Upon get config, it returns pointers to weight matrices along with their shape information in same structure.
XA_NNLIB_LSTM_BIAS	<code>xa_nnl-lib_lstm_biases_t</code>	RW	NA	NA	Bias vectors, pointers to bias vectors along with shape information needs to be passed via <code>xa_nnl-lib_lstm_biases_t</code> structure for set config. Upon get config, it returns pointers to bias vectors along with their shape information in same structure.
XA_NNLIB_LSTM_INPUT_SHAPE	<code>xa_nnl-lib_shape_t</code>	R	NA	NA	Input shape information, get information of the input shape expected by the layer.
A_NNLIB_LSTM_OUTPUT_SHAPE	<code>xa_nnl-lib_shape_t</code>	R	NA	NA	Output shape information, get information of the output shape expected by layer.

4.3 CNN Layer

The CNN APIs are defined in `xa_nnlb_cnn_api.h`.

4.3.1 CNN Layer Specification

The CNN layer implements Standard 2D Convolution, Standard 1D Convolution, and Depthwise Separable 2D Convolution. Refer to the equations in Section 3.2.1 for Standard 2D Convolution, Section 3.2.2 for Standard 1D Convolution, and Section 3.2.4 for Depthwise Separable 2D Convolution.

4.3.2 Error Codes Specific to CNN

Other than common error codes explained in Section 2.3, the CNN layer may also report the following error codes, which may be generated during the initialization stage.

- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_ALGO`
Algorithm is not supported
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_PRECISION`
I/O precision is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_BIAS_SHIFT`
Value of Bias shift is not supported
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_ACC_SHIFT`
Value of Accumulator shift is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_STRIDE`
Value of strides is not supported
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_PADDING`
Value of padding is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_INPUT_SHAPE`
Input shape dimension is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_OUTPUT_SHAPE`
Out shape dimension is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_KERNEL_SHAPE`
Kernel shape dimension is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_BIAS_SHAPE`
Bias shape dimension is not supported.
- `XA_NNLB_CNN_CONFIG_FATAL_INVALID_PARAM_ID`
Parameter identifier (`param_id`) is not valid

- **XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_COMBINATION**

Parameter combination (param_id) is not valid

The following error codes may be generated during the execution stage.

- **XA_NNLIB_CNN_CONFIG_FATAL_INVALID_INPUT_SHAPE**

Input shape passed during execution does not match with the input shape passed during initialization

4.3.3 API Functions Specific to CNN

4.3.3.1 Query Functions

Table 4-23 CNN Get Persistent Size Function

Function	<code>xa_nnl-lib-cnn-get-persistent-fast</code>
Syntax	<pre>Int32 xa_nnl-lib-cnn-get-persistent-fast (xa_nnl-lib-cnn-init-config_t *config)</pre>
Description	Returns persistent memory size in bytes required by CNN layer.
Parameters	Input: <code>config</code> Initial configuration parameters (see Table 4-29).
Errors	If return value is less than 0, then it is an error. Following are the possible error codes: <ul style="list-style-type: none"> ■ XA_NNLIB_FATAL_MEM_ALLOC ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_ALGO Algorithm is not supported ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PRECISION I/O precision is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_BIAS_SHIFT Value of Bias shift is not supported ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_ACC_SHIFT Value of Accumulator shift is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_STRIDE Value of strides is not supported ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PADDING Value of padding is not supported.

	<ul style="list-style-type: none"> ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_INPUT_SHAPE Input shape dimension is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_OUTPUT_SHAPE Out shape dimension is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_KERNEL_SHAPE Kernel shape dimension is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_BIAS_SHAPE Bias shape dimension is not supported ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_ID Parameter identifier (param_id) is not valid ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_COMBINATION Parameter combination (param_id) is not valid
--	--

Table 4-24 CNN Get Scratch Size Function

Function	<code>xa_nnlb_cnn_get_scratch_fast</code>
Syntax	<pre>Int32 xa_nnlb_cnn_get_scratch_fast (xa_nnlb_cnn_init_config_t *config)</pre>
Description	Returns scratch memory size in bytes required by CNN layer.
Parameters	Input: <code>config</code> Initial configuration parameters (see Table 4-29).
Errors	If return value is less than 0, then it is an error. Following are the possible error codes: <ul style="list-style-type: none"> ■ XA_NNLIB_FATAL_MEM_ALLOC ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_ALGO Algorithm is not supported ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PRECISION I/O precision is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_BIAS_SHIFT Value of bias shift is not supported ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_ACC_SHIFT Value of Accumulator shift is not supported. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_STRIDE Value of strides is not supported

- | | |
|--|--|
| | <ul style="list-style-type: none">■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PADDING</code>
Value of padding is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_INPUT_SHAPE</code>
Input shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_OUTPUT_SHAPE</code>
Out shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_KERNEL_SHAPE</code>
Kernel shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_BIAS_SHAPE</code>
Bias shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_ID</code>
Parameter identifier (param_id) is not valid■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_COMBINATION</code>
Parameter combination (param_id) is not valid |
|--|--|

4.3.3.2 Initialization Stage

Table 4-25 CNN Init Function

Function	<code>xa_nnlib_cnn_init</code>
Syntax	<pre>int xa_nnlib_cnn_init (xa_nnlib_handle_t handle, xa_nnlib_cnn_init_config_t *config)</pre>
Description	Reset the CNN layer API handle into its initial state. Set up the CNN layer to the specified initial configuration parameters.
Parameters	<p>Input: <code>handle</code> Pointer to the component persistent memory. This is the opaque handle. Required size: see <code>xa_nnlib_cnn_get_persistent_fast</code>. Required alignment: 8 bytes.</p> <p>Input: <code>config</code> Initial configuration parameters (see Table 4-29). Note that the initial configuration parameters <i>must</i> be identical to those passed to query functions.</p>
Errors	<p>If the return value is not <code>XA>NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA>NNLIB_FATAL_MEM_ALLOC</code> One of the pointers is invalid. ■ <code>XA>NNLIB_FATAL_MEM_ALIGN</code> One of the pointers is not properly aligned. ■ <code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_ALGO</code> Algorithm is not supported. ■ <code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_PRECISION</code> I/O precision is not supported. ■ <code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_BIAS_SHIFT</code> Value of Bias shift is not supported. ■ <code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_ACC_SHIFT</code> Value of Accumulator shift is not supported. ■ <code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_STRIDE</code> Value of strides is not supported. ■ <code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_PADDING</code> Value of padding is not supported.

- | | |
|--|--|
| | <ul style="list-style-type: none">■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_INPUT_SHAPE</code>
Input shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_OUTPUT_SHAPE</code>
Out shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_KERNEL_SHAPE</code>
Kernel shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_BIAS_SHAPE</code>
Bias shape dimension is not supported.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_ID</code>
Parameter identifier (<code>param_id</code>) is not valid.■ <code>XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_COMBINATION</code>
Parameter combination (<code>param_id</code>) is not valid. |
|--|--|

4.3.3.3 Execution Stage

Table 4-26 CNN Execution Function

Function	<code>xa_nnlib_cnn_process</code>
Syntax	<pre>int xa_nnlib_cnn_process (xa_nnlib_handle_t handle, void *scratch, void *input, void *output, xa_nnlib_shape_t *p_in_shape, xa_nnlib_shape_t *p_out_shape)</pre>
Description	Processes one input shape to generate one output shape.
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>scratch</code> A pointer to the scratch buffer. Required alignment: 8 bytes.</p> <p>Input: <code>input</code> A pointer to the input buffer. Input buffer contains input data. Required alignment: 8 bytes.</p> <p>Output: <code>output</code> A pointer to the output buffer. Output is written to the output buffer. Required alignment: 8 bytes.</p> <p>Input/Output: <code>p_in_shape</code> Pointer to the shape containing input buffer dimensions. Contains the length of input data passed to the CNN layer. Required alignment: 4 bytes.</p> <p>Output: <code>p_out_shape</code> Pointer to the shape for output buffer dimensions. Upon return, <code>*p_out_shape</code> is filled with the length of output generated by the CNN layer. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA>NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA>NNLIB_FATAL_MEM_ALLOC</code> One of the pointers is NULL

	<ul style="list-style-type: none"> ■ XA_NNLIB_FATAL_MEM_ALIGN One of the pointers is not having required alignment ■ XA_NNLIB_FATAL_INVALID_SHAPE Input shape passed during execution does not match with the input shape passed during initialization
--	--

Table 4-27 CNN Set Parameter Function Details

Function	<code>xa_nnl-lib_cnn_set_config</code>
Syntax	<pre>int xa_nnl-lib_cnn_set_config (xa_nnl-lib_handle_t handle, xa_nnl-lib_cnn_param_id_t param_id, void *params)</pre>
Description	Sets the parameter specified by <code>param_id</code> to the value passed in the buffer pointed to by <code>params</code> .
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>param_id</code> Identifies the parameter to be written. Refer to Table 4-32 for the list of supported parameters.</p> <p>Input: <code>params</code> A pointer to a buffer that contains the parameter value. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA_NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ XA_NNLIB_FATAL_MEM_ALLOC One of the pointers (<code>handle</code> or <code>params</code>) is <code>NULL</code>. ■ XA_NNLIB_FATAL_MEM_ALIGN One of the pointers (<code>handle</code> or <code>params</code>) is not aligned correctly. ■ XA_NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_ID Parameter identifier (<code>param_id</code>) is not valid.

Table 4-28 CNN Get Parameter Function Details

Function	<code>xa_nnlib_cnn_get_config</code>
Syntax	<pre>int xa_nnlib_cnn_get_config(xa_nnlib_handle_t handle, xa_nnlib_cnn_param_id_t param_id, void *params)</pre>
Description	Gets the value of the parameter specified by <code>param_id</code> in the buffer pointed to by <code>params</code> .
Parameters	<p>Input: <code>handle</code> The opaque component handle. Required alignment: 8 bytes.</p> <p>Input: <code>param_id</code> Identifies the parameter to be read. Refer to Table 4-32 for the list of supported parameters.</p> <p>Output: <code>params</code> A pointer to a buffer that is filled with the parameter value when the function returns. Required alignment: 4 bytes.</p>
Errors	<p>If the return value is not <code>XA>NNLIB_NO_ERROR</code>, it implies that the function has encountered one of the following errors:</p> <ul style="list-style-type: none"> ■ <code>XA>NNLIB_FATAL_MEM_ALLOC</code> One of the pointers (<code>handle</code> or <code>params</code>) is <code>NULL</code>. ■ <code>XA>NNLIB_FATAL_MEM_ALIGN</code> One of the pointers (<code>handle</code> or <code>params</code>) is not aligned correctly. <p><code>XA>NNLIB_CNN_CONFIG_FATAL_INVALID_PARAM_ID</code> Parameter identifier (<code>param_id</code>) is not valid.</p>

4.3.4 Structures Specific to CNN

Table 4-29 CNN Config Structure xa_nnlib_cnn_init_config_t

Element Type	Element Name	Range	Default	Description
xa_nnlib_shape_t	input_shape	NA	height = 16 width = 16 channels = 4	Input shape dimensions
Int32	output_height	NA	16	Output height
Int32	output_width	NA	16	Output width
Int32	output_channels	NA	4	Output depth or channels
Int32	output_format	0 or 1	0	Output data format 0: SHAPE_CUBE_DWH_T 1: SHAPE_CUBE_WHD_T
xa_nnlib_shape_t	kernel_std_shape	NA	height = 16 width = 16 channels = 4	Standard 1D/2D Convolution Kernel (Filter) shape dimensions output_channels indicate number of kernels
xa_nnlib_shape_t	kernel_ds_depth_shape	NA	NA	Depthwise Separable 2D Convolution - Depthwise Kernel (filter) Dimensions
xa_nnlib_shape_t	kernel_ds_point_shape	NA	NA	Depthwise Separable 2D Convolution - Pointwise Kernel (filter) Dimensions
xa_nnlib_shape_t	bias_std_shape	NA	channels = 4	Standard 1D/2D Convolution Bias dimensions
xa_nnlib_shape_t	bias_ds_depth_shape	NA	NA	Depthwise Separable 2D Convolution - Depthwise Bias) Dimensions
xa_nnlib_shape_t	bias_ds_point_shape	NA	NA	Depthwise Separable 2D Convolution – Pointwise Bias Dimensions
xa_nnlib_cnn_precision_t	precision	XA_NNLIB_CNN_16bx16b, XA_NNLIB_CNN_8bx16b, XA_NNLIB_CNN_8bx8b, XA_NNLIB_CNN_f32xf32	XA_NNLIB_CNN_8bx16b	Kernel (filter), input, output precision setting
Int32	bias_shift	-31 to 31	7	Q-format adjustment for bias before addition into accumulator, +/- value - left/right shift

Element Type	Element Name	Range	Default	Description
Int32	acc_shift	-31 to 31	-7	Q-format adjustment for accumulator before rounding to result, +/- value - left/right shift
Int32	channels_multiplier	NA	NA	Depthwise Separable 2D Convolution - channel multiplier. (channels_multiplier * input_channels) must be multiple of 4
Int32	x_padding	NA	2	Left side padding to be added to input
Int32	y_padding	NA	2	Top padding to be added to input
Int32	x_stride	NA	2	Strides over padded input in width dimension
Int32	y_stride	NA	2	Strides over padded input in height dimension
xa_nnlb_cnn_algo_t	algo	NA	XA_NNLB_CNN_CO NV2D_STD	Convolution algorithm

4.3.5 Enums Specific to CNN

Table 4-30 Enum xa_nnlb_cnn_precision_t

Element	Description
XA_NNLB_CNN_16b×16b	Coef: 16 bits, I/O: 16 bits fixed point
XA_NNLB_CNN_8b×16b	Coef: 8 bits, I/O: 16 bits fixed point
XA_NNLB_CNN_8b×8b	Coef: 8 bits, I/O: 8 bits fixed point
XA_NNLB_CNN_f32×f32	Coef: single precision float, I/O: single precision float

Table 4-31 Enum xa_nnlb_cnn_algo_t

Element	Description
XA_NNLB_CNN_CONV1D_ST	Standard 1D Convolution
XA_NNLB_CNN_CONV2D_STD	Standard 2D Convolution
XA_NNLB_CNN_CONV2D_DS	Depthwise Separable 2D Convolution

Table 4-32 describes parameter IDs for parameters supported by CNN. It contains the following columns:

- Parameter ID: Parameter identifier (`param_id`).
- Value type: A pointer (`params`) to a variable of this type is to be passed.
- RW: Indicates whether the parameter can be read (`get`) and/or written (`set`).
- Range: Indicates valid values of the parameter.
- Default: Default value of the parameter
- Description: Brief description of the parameter.

Table 4-32 CNN Specific Parameters

Parameter ID	Value Type	RW	Range	Default	Description
XA>NNLIB_CNN_KERNEL	<code>vect_t</code> []	<i>RW</i>	NA	NA	Kernel shape information, get or set information of the kernel shape expected by the layer
XA>NNLIB_CNN_BIAS	<code>vect_t</code> []	<i>RW</i>	NA	NA	Bias shape information, get or set information of the bias shape expected by the layer
XA>NNLIB_CNN_INPUT_SHAPE	<code>xa_nnl-lib_shape_t</code>	<i>R</i>	NA	NA	Input shape information, get information of the input shape expected by the layer.
XA>NNLIB_CNN_OUTPUT_SHAPE	<code>xa_nnl-lib_shape_t</code>	<i>R</i>	NA	NA	Output shape information, get information of the output shape produced by layer.

5. Introduction to the Example Testbench

5.1 Making the Library

If you have source code distribution, you must build the NN library before you can build the testbench. To do so, follow these steps:

1. Go to `build`.
2. From the command prompt, enter:

```
xt-make -f makefile clean all install
```

The NN library `xa_nnlib.a` will be built and copied to the `lib` directory.

5.2 Making the Executable

To build the testbenches, follow these steps:

1. Go to `test/build`.
2. From the command-line prompt, enter:

```
xt-make -f makefile_testbench_sample clean all
```

This will build the example testbenches for all the kernels and layers.

The following header files are common and used by all testbenches.

- Testbench header files (`test/include`)
 - `xt_profiler.h`
 - `cmdline_parser.h`
 - `file_io.h`
 - `xt_manage_buffers.h`

The following sections describe each low-level kernel and layer testbench.

5.3 Sample Testbench for Matrix X Vector Multiplication Kernels

The NN library Matrix X Vector Multiplication Kernels are provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_matXvec_testbench.c

5.3.1 Usage

The NN library Matrix X Vector Multiplication Kernels executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_matXvec_test [options]
```

Following are available options:

Option	Description	Additional Information
-rows	Rows of mat1.	
-cols1	Columns of mat1 and rows of mat2 (Default=32)	Columns of mat1 must be multiple of 4
-cols2	Columns of mat2 (Default=32)	Columns of mat2 must be multiple of 4
-row_stride1	Row stride for mat1(Default=32)	
-row_stride2	Row stride for mat2(Default=32)	
-vec_count	Vec count for Time batching(Default=1)	
-acc_shift	Accumulator shift(Default=-7)	
-bias_shift	Bias shift(Default=7)	
-mat_precision	8, 16, -1(single precision float), -3 (asym8u) or -5 (sym8s); (Default=16)	
-inp_precision	8, 16, -1(single precision float), -3 (asym8u) or -4 (asym8s); (Default=16)	
-out_precision	8, 16, 32, 64 -1(single precision float), -3 (asym8u) or -4 (asym8s); (Default=16)	
-bias_precision	16, 64 -1(single precision float), 32(for quantized 8-bit kernels); (Default=16)	
-membank_padding	0 or 1 (Default=1)	
-frames	Positive number; (Default=2)	
-activation	Sigmoid, tanh, relu or softmax (Default= bypass i.e. no activation for output)	

Option	Description	Additional Information
-write_file	Set to 1 to write input and output vectors to file; (Default=0)	
-read_inp_file_name	Full filename for reading inputs (order - mat1, vec1, mat2, vec2, bias)	
-read_ref_file_name	Full filename for reading reference output	
-write_inp_file_name	Full filename for writing inputs (order - mat1, vec1, mat2, vec2, bias)	
-write_out_file_name	Full filename for writing output	
-verify	Verify output against provided reference	0: Disable, 1: Bit exact match (Default=1)
-batch	Flag to execute time batching kernels	0: Disable, 1: Enable (Default=0)
-fc	Flag to execute fully connected kernels	0: Disable, 1: Enable (Default=0)
-help	Prints help	

If no command line arguments are given, the Matrix X Vector Multiplication Kernels sample testbench runs with default values.

5.4 Sample Testbench for Convolution Kernels

The NN library convolutional kernels are provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_conv_testbench.c

5.4.1 Usage

The NN Library convolutional kernels executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_conv_test [options]
```

Following are available options:

Option	Description
-input_height	Input height (Default=16)
-input_width	Input width (Default=16)
-input_channels	Input channels (Default=4)
-kernel_height	Kernel height (Default=3)

Option	Description
-kernel_width	Kernel width (Default=3)
-out_channels	Out channels (Default=4)
-channels_multiplier	Channel Multiplier (Default=1)
-x_stride	Stride in width dimension (Default=2)
-y_stride	Stride in height dimension (Default=2)
-x_padding	Left padding in width dimension (Default=2)
-y_padding	Top padding in height dimension (Default=2)
-dilation_height	Kernel Dilation factor in height dimension (Default=1)
-dilation_width	Kernel Dilation factor in width dimension (Default=1)
-out_height	Output height (Default=16)
-out_width	Output width (Default=16)
-bias_shift	Bias shift (Default=7)
-acc_shift	Accumulator shift (Default=-7)
-inp_data_format	0 (SHAPE_CUBE_DWH_T), 1 (SHAPE_CUBE_WHD_T) (Default=1)
-out_data_format	0 (SHAPE_CUBE_DWH_T), 1 (SHAPE_CUBE_WHD_T) (Default=0)
-inp_precision	8, 16, -1(single precision float), -3 (asym8u) or -4 (asym8s); (Default=16)
-kernel_precision	8, 16, -1(single precision float), -3 (asym8u) or -5 (sym8s); (Default=8)
-out_precision	8, 16, -1(single precision float), -3 (asym8u) or -4 (asym8s); (Default=16)
-bias_precision	8, 16, -1(single precision float), 32(for quantized 8-bit kernels) (Default=16)
-frames	Positive number (Default=2)
-kernel_name	conv2d_std, conv2d_depth, conv1d_std, dilated_conv2d_std; (Default= conv2d_std)
-write_file	Set to 1 to write input and output vectors to file; (Default=0)
-read_inp_file_name	Full filename for reading inputs (order - input, kernel, bias, (pointwise kernel, pointwise bias for depth separable))

Option	Description
-read_ref_file_name	Full filename for reading reference output
-write_inp_file_name	Full filename for writing inputs (order - input, kernel, bias, (pointwise kernel, pointwise bias for depth separable))
-write_out_file_name	Full filename for writing output
-verify	Verify output against provided reference; 0: Disable, 1: Bit exact match (Default=1)
-help	Prints help

If no command line arguments are given, the Convolutional Kernels sample testbench runs with default values.

5.5 Sample Testbench for Activation Kernels

The NN library activation kernels are provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_activations_testbench.c

5.5.1 Usage

The NN library activation kernels executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_activation_test [options]
```

Following are available options:

Option	Description
-num_elements	Number of elements (Default=32)
-relu_threshold	Threshold for relu in Q16.15 (Default= 32768 i.e. =1 in Q16.15)
-inp_precision	8, 16, 32, -1 (single precision float), -3 (asym8u) or -4 (asym8s); (Default=32)
-out_precision	8, 16, 32, -1(single precision float), -3 (asym8u) or -4 (asym8s); (Default=32)
-frames	Positive number (Default=2)
-activation	Sigmoid, tanh, relu, relu_std, relu1, relu6, activation_min_max, softmax, hard_swish or prelu (Default= sigmoid)
-write_file	Set to 1 to write input and output vectors to file; (Default=0)
-read_inp_file_name	Full filename for reading input
-read_ref_file_name	Full filename for reading reference output
-write_inp_file_name	Full filename for writing input
-write_out_file_name	Full filename for writing output
-verify	Verify output against provided reference; 0: Disable, 1: Bit exact match (Default=1)
Quantized 8-bit specific parameters	
-diffmin	Diffmin; Default=-15
-input_left_shift	Input_left_shift; Default=27
-input_multiplier	Input_multiplier; Default=2060158080

Option	Description
-activation_max	asym8u input data activation max; Default=0
-activation_min	asym8u input data activation min; Default=0
-input_range_radius	Sigmoid_asym8u input parameter; Default=128
-zero_point	Sigmoid_asym8u input parameter; Default=0
-inp_zero_bias	Zero bias value for input Default=0
-alpha_zero_bias	Prelu parameter - Zero bias value for alpha Default=0
-alpha_multiplier	Prelu parameter - Multiplier value for alpha Default=0x40000000
-alpha_shift	Prelu parameter - Shift value for alpha Default=0
-reluish_multiplier	Hard Swish parameter - Multiplier value for relu scale Default=0x40000000
-reluish_shift	Hard Swish parameter - Shift value for relu scale Default=0
-out_multiplier	Multiplier value for output Default=0x40000000
-out_shift	Shift value for output Default=0
-out_zero_bias	Zero bias value for output Default=0
-help	Prints help

If no command line arguments are given, the Activation Kernels sample testbench runs with default values.

5.6 Sample Testbench for Pooling Kernels

The NN library pooling kernels are provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_pool_testbench.c

5.6.1 Usage

The NN library pooling kernels executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_pool_test [options]
```

Following are available options:

Option	Description
-input_height	Input height (Default=16)
-input_width	Input width (Default=16)
-input_channels	Input channels (Default=4)
-kernel_height	Kernel height (Default=3)
-kernel_width	Kernel width (Default=3)
-x_stride	Stride in width dimension (Default=2)
-y_stride	Stride in height dimension (Default=2)
-x_padding	Left padding in width dimension (Default=2)
-y_padding	Top padding in height dimension (Default=2)
-out_height	Output height (Default=16)
-out_width	Output width (Default=16)
-acc_shift	Accumulator shift (Default=-7)
-inp_data_format	0 (SHAPE_CUBE_DWH_T), 1 (SHAPE_CUBE_WHD_T) (Default=1)
-out_data_format	0 (SHAPE_CUBE_DWH_T), 1 (SHAPE_CUBE_WHD_T) (Default=1)
-inp_precision	8, 16, -1(single precision float), -3(asym8u) (Default=16)
-out_precision	8, 16, -1(single precision float), -3(asym8u) (Default=16)
-frames	Positive number (Default=2)
-kernel_name	avgpool, maxpool (Default=avgpool)
-write_file	set to 1 to write input and output vectors to file; (Default=0)
-read_inp_file_name	Full filename for reading inputs (order - inp)
-read_ref_file_name	Full filename for reading reference output
-write_inp_file_name	Full filename for writing inputs (order - inp)
-write_out_file_name	Full filename for writing output
-verify	Verify output against provided reference; 0: Disable, 1: Bit exact match (Default=1)
-help	Prints help

If no command line arguments are given, the Pooling Kernels sample testbench runs with default values.

5.7 Sample Testbench for Basic Operations Kernels

The NN library basic kernels are provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_basic_testbench.c

5.7.1 Usage

The NN library basic kernels executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_basic_test [options]
```

Following are available options:

Option	Description
-io_length	Input/output vector length; Default=1024
-inp_precision	16, -3 (asym8u), -1 (single prec float), -4(asym8s), 1(bool); Default=-1
-out_precision	-3 (asym8u), -1 (single prec float), -4(asym8s) , 1(bool); Default=-1
-vec_count	Number of input vectors; Default=1
-frames	Positive number; Default=2
-kernel_name	elm_add, elm_sub, elm_mul, elm_floor, dot_prod, elm_min and elm_max, elm_equal, elm_notequal, elm_greater, elm_greaterequal, elm_less, elm_lessequal, elm_logicaland, elm_logicalor, elm_logicalnot, reduce_max_4D, reduce_mean_4D, elm_min_4D_Bcast, elm_max_4D_Bcast; Default=elem_add
-write_file	Set to 1 to write input and output vectors to file; Default=0
-read_inp1_file_name	Full filename for reading inputs (order - inp)
-read_inp2_file_name	Full filename for reading inputs (order - inp)
-read_ref_file_name	Full filename for reading reference output
-write_inp1_file_name	Full filename for writing inputs (order - inp)

Option	Description
-write_inp2_file_name	Full filename for writing inputs (order - inp)
-write_out_file_name	Full filename for writing output
-verify	Verify output against provided reference; 0: Disable, 1: Bit exact match; Default=1
-read_inp_shape_str	Takes the input shape dimensions(space ' ' separated) as a string
-read_out_shape_str	Takes the output shape dimensions(space ' ' separated) as a string
-read_axis_data_str	Takes the axis data (space ' ' separated) as a string
Broadcast specific parameters	
-input1_numElements	Number of elements in input
-input2_numElements	Number of elements in input
-input1_strides	Input strides
-input2_strides	Input strides
Quantized 8-bit specific parameters	
-output_zero_bias	Output zero bias; Default=127
-output_left_shift	Output_left_shift; Default=1
-output_multiplier	Output_multiplier; Default=0x7fff
-output_activation_min	Output_activation_min; Default=0
-output_activation_max	Output_activation_max; Default = 225
-input1_zero_bias	Input1 zero bias; Default=-127
-input1_left_shift	Input1 left shift; Default=0
-input1_multiplier	Input1 multiplier; Default=0x7fff
-input2_zero_bias	Input2 zero bias; Default=-127
-input2_left_shift	Input2 left shift; Default=0
-input2_multiplier	Input2 multiplier; Default=0x7fff
-left_shift	Global left shift; Default=0
-h	Prints help

If no command line arguments are given, the Basic Kernels sample testbench runs with default values.

5.8 Sample Testbench for Normalization Kernels

The NN library Normalization Kernels are provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_norm_testbench.c

5.8.1 Usage

The NN library Normalization Kernels executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_norm_test [options]
```

Following are available options:

Option	Description
-num_elms	Number of elements; Default=256
-inp_precision	-4(asym8s); Default=16
-out_precision	-4(asym8s); Default=16
-frames	Positive number; Default=2
-kernel_name	L2_norm; Default=l2_norm
-zero_point	Input Zero point; Default = 0
-write_file	Set to 1 to write input and output vectors to file; Default=0
-read_inp_file_name	Full filename for reading inputs (order - inp)
-read_ref_file_name	Full filename for reading reference output
-write_inp_file_name	Full filename for writing inputs (order - inp)
-write_out_file_name	Full filename for writing output
-verify	Verify output against provided reference; 0: Disable, 1: Bit exact match; Default=1
-h	Prints help

If no command line arguments are given, the Normalization Kernels sample testbench runs with default values.

5.9 Sample Testbench for GRU Layer

The NN library GRU layer is provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)
 - xa_nn_gru_testbench.c

5.9.1 Usage

The NN library GRU executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_gru_test [options]
```

Following are available options:

Option	Description	Additional Information
--in_feats	Input length (Default=256)	Range: 4-2048 Note: Input length must be multiple of 4
--out_feats	Output length (Default=256)	Range: 4-2048 Note: Output length must be multiple of 4
--membank_padding	Memory bank padding (Default=1)	Must be 0 or 1
--mat_prec	Coefficient precision (Default=16)	Must be 8 or 16
--vec_prec	Input precision (Default=16)	Must be 16
--verify	Verify output against ref output (Default=1)	Supported values: 0: Disable, 1: Enable
--input_file	Input file name	
--filter_path	Path where file containing filter are stored	
--output_file	File to which output will be written	
--prev_h_file	File containing context data	
--ref_file	File which has ref output	
-help	Prints help	

If no command line arguments are given, the GRU sample testbench runs with default values.

5.10 Sample Testbench for LSTM Layer

The NN library LSTM layer is provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (test/src)

- `xa_nn_lstm_testbench.c`

5.10.1 Usage

The NN library LSTM executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_lstm_test [options]
```

Following are available options:

Option	Description	Additional Information
<code>--in_feats</code>	Input length (Default=256)	Range: 4-2048 Note: Input length must be multiple of 4
<code>--out_feats</code>	Output length (Default=256)	Range: 4-2048 Note: Output length must be multiple of 4
<code>--membank_padding</code>	Memory bank padding (Default=1)	Must be 0 or 1
<code>--mat_prec</code>	Coefficient precision (Default=16)	Must be 8 or 16
<code>--vec_prec</code>	Input precision (Default=16)	Must be 16
<code>--verify</code>	Verify output against ref output (Default=1)	Supported values: 0: Disable, 1: Enable
<code>--input_file</code>	File containing input shape	
<code>--filter_path</code>	Path where file containing filter are stored	
<code>--output_file</code>	File to which output will be written	
<code>--output_cell_file</code>	File to which cell output will be written	
<code>--prev_h_file</code>	File containing context (previous output) data	
<code>--prev_c_file</code>	File containing context (previous cell state) data	
<code>--ref_file</code>	File which has ref output	
<code>--ref_cell_file</code>	File which has ref cell output	
<code>-help</code>	Prints help	

If no command line arguments are given, the LSTM sample testbench runs with default values.

5.11 Sample Testbench for CNN Layer

The NN library CNN layer is provided with a sample testbench application. The supplied testbench consists of the following files:

- Testbench source files (`test/src`)
 - `xa_nn_cnn_testbench.c`

5.11.1 Usage

The NN Library CNN executable can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_cnn_test [options]
```

Following are available options:

Option	Description
-input_height	Input height (Default=16)
-input_width	Input width (Default=16)
-input_channels	Input channels (Default=4)
-kernel_height	Kernel height (Default=3)
-kernel_width	Kernel width (Default=3)
-out_channels	Out channels (Default=4)
-channels_multiplier	Channel Multiplier (Default=1)
-x_stride	Stride in width dimension (Default=2)
-y_stride	Stride in height dimension (Default=2)
-x_padding	Left padding in width dimension (Default=2)
-y_padding	Top padding in height dimension (Default=2)
-out_height	Output height (Default=16)
-out_width	Output width (Default=16)
-bias_shift	Bias shift (Default=7)
-acc_shift	Accumulator shift (Default=-7)
-out_data_format	Output data format, 0 (SHAPE_CUBE_DWH_T), 1 (SHAPE_CUBE_WHD_T); (Default=0)
-inp_precision	8, 16, -1(single precision float); (Default=16)
-kernel_precision	8, 16, -1(single precision float); (Default=8)
-out_precision	8, 16, -1(single precision float); (Default=16)
-bias_precision	8, 16, -1(single precision float); (Default=16)
-frames	Positive number; (Default=2)
-kernel_name	conv2d_std, conv2d_depth, conv1d_std; (Default=conv2d_std)
-write_file	Set to 1 to write input and output vectors to file; (Default=0)

Option	Description
-read_inp_file_name	Full filename for reading inputs (order - input, kernel, bias, (pointwise kernel, pointwise bias for depth separable))
-read_ref_file_name	Full filename for reading reference output
-write_inp_file_name	Full filename for writing inputs (order - input, kernel, bias, (pointwise kernel, pointwise bias for depth separable))
-write_out_file_name	Full filename for writing output
-verify	Verify output against provided reference; 0: Disable, 1: Bit exact match; Default=1
-help	Prints help

If no command line arguments are given, the CNN sample testbench runs with default values.

5.12 Sample Testbench of TF Micro Lite Example

The NN library package contains a sample testbench of the TF Micro Lite example application. This application is based on the Micro Speech Example provided by Google [2] in the TensorFlow GitHub repository (TensorFlow v2.3). It accepts 1 sec wav audio from the input file and outputs the scores for the twelve classes: silence, unknown, yes, no, up, down, left, right, on, off, stop, go. It uses optimized NN Library kernels for the inference as per the network model specified in the application. The TensorFlow Lite model 'low_latency_conv', comprising of one depthwise convolutional layer followed by three fully connected layers is used in this application. The model uses three operations depthwiseConv, FullyConnected and Softmax. The testbench consists of the following files:

- Testbench source files (test/tf_micro_lite)
 - test/src/... The TF Micro Lite application source
 - tensorflow/... Supporting files from the tensorflow package

5.12.1 Usage

The TF Micro Lite example can be run with command-line options as follows.

```
$ xt-run [--mem_model] [--turbo] xa_nn_tfulite_test <input-wav-file.wav>
<test_result.txt>
```

Option	Description
<input-wav-file.wav>	Full filename for reading input wav audio
<test_result.txt>	Output text file for storing the classification score
-help	Shows help

If no command line arguments are given, the sample testbench runs with default values.

6. References

- [1] Reference Wiki page for GRU. https://en.wikipedia.org/wiki/Gated_recurrent_unit

- [2] TF Micro Lite speech recognition example:
https://github.com/tensorflow/tensorflow/tree/r2.3/tensorflow/lite/micro/examples/micro_speech

- [3] [TensorFlow Lite for Microcontrollers](#)

- [4] TensorFlow XLA Documentation: <https://www.tensorflow.org/xla/broadcasting>
NumPy Theory: <https://numpy.org/devdocs/user/theory.broadcasting.html>
General Broadcasting syntax: <https://www.tensorflow.org/guide/tensor#broadcasting>

- [5] 'strides' as defined in the structure 'NDArrayDesc' at
<https://github.com/tensorflow/tensorflow/blob/master/tensorflow/lite/kernels/internal/common.h>