

Yellow Of The Egg
Lukas Baischer
Add your names ..

SoC Design Laboratoy
384.157, Winter Term 2019

MNIST-FPGA Specification

Maybe add an additional Input Layer which is responsible to communicate with the DMA and converts the data from 32 bit to 8 bit and sends it to the memory controller

- Create a bitstream with the generated vhdl code

Requirements of the Trainings Software:

- Sends image data to Zedboard
- Receives results from Zedboard
- Create a figure of accuracy and performance
- Optional: Send bitstream to hardware which updates the bitstream

3.1.1 Interface to Zedboard

Ethernet is used for the communication of the remote host system and the embedded Linux which is running on the Zedboard.

The embedded Linux distribution running on the board should automatically receive an IP address when connected to a network. When in doubt the address can be found out with the `ifconfig` command.

The software has a client-server model with the embedded system acting as a server and the host as a client. Once running, the server software is listening for new outside connections.

Different types of data need to be transmitted:

- The 28x28 input images showing digits between 0 and 9 is transferred from host to Zedboard.
- The probability of resulting numbers between 0 and 9 is transmitted from Zedboard to host.
- control and status signals in both directions
- Optional: Bitstream file for dynamically update the bitstream at the Zedboard

Who is the host and client now

3.1.2 Notes

On Windows host systems, *Network Discovery* needs to be enabled and in some cases a Firewall exception for the used ports needs to be set for a connection to be established.

3.2 ARM Top-Level software

The ARM top-level software receives the image data from a remote device and sends the results back to this device. Control of the hardware.

Optional feature: Update Bitstream file using `/dev/xdevcfg`

Requirements of the ARM Top-Level Software:

- Receive image data
- send results to remote PC
- Send and receive control signals from remote PC
- Send image data to driver user layer and receive results from driver user layer
- Send and receive status and control signals to driver user layer
- Run at start-up

Add more information and specify the requirements

3.2.1 Interface to remote PC

See Section 3.1.1.

3.2.2 Interface to kernel layer

Python wrapper are used for the interface between the top level software which is programmed in python and the hardware drivers which are programmed in C

Add more information and specify the requirements of the interface

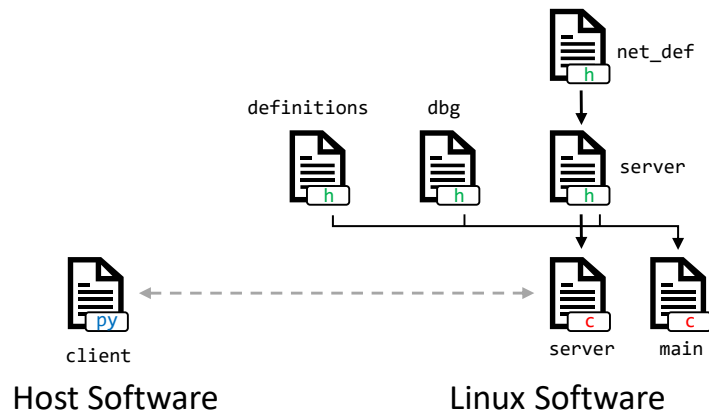


Figure 2: File tree for the software

3.2.3 File Tree of ARM Top-Level software

- `net_def.h` Contains definitions for networking, e.g. ports used.
- `dbg.h` Contains debugging macros for logging and error handling.
- `definitions.h` Contains information about the neural network, e.g. the number and type of Convolutional Neural Network (CNN) stages, layers in the fully connected network, input size and so on.
- `server.{c,h}` Handles the connection with the host software.
- `main.c` Contains the `main()` function with the main program loop that transmits and manages data to the hardware and from the host system.
- `client.py` Handles the connection with the client software.

Update this section. Do we still use the C code or do we plan to implement everything in python?

3.3 User Layer Driver Software

The user layer driver software implements an interface between the ARM Top-Level software and the driver for the programmable logic. It is implemented in C. It is supposed to handle the entire communication with the driver so that the hardware is only abstractly visible for the ARM Top-Level software.

For example the ARM top-level software sees the network as a class in python which has a method `load_new_image` data with a numpy array as input and a finish signal as a output. This method should call the user layer driver software which handles the communication between user space and kernel space. In a similar way each IP should be a class in python.

Requirements of the User Layer Driver Software:

- Communication with the kernel space drivers
- Use python wrapper to communicate with ARM Top-Level software
- Easy to use interface from Top-Level
- No knowledge of the hardware should be necessary to use the interface
- Data encapsulation to avoid the Top-Level Software from corrupting the memory

3.3.1 File Tree of User Layer Driver Software

4 Hardware

4.1 Memory Controller

The task of the memory controller is to provide valid data for the NN-layers. It communicates with the Block-Ram. The memory controller is responsible for ensuring that the next layer has valid data at all times. The second task of the memory controller is to save the data of the previous data in a free memory address in the Block-RAM.

Would be nice if we have something similar as in 3.2.3

Is it better to have the shiftregister, we discussed last time in

4.1.1 Interfaces

- S_LAYER: interface to previous layer

signal	direction	type	width	description
--------	-----------	------	-------	-------------

- M_LAYER: interface to next layer

signal	direction	type	width	description
--------	-----------	------	-------	-------------

- BRAM_PORTA: write interface to BRAM

signal	direction	type	width	description
--------	-----------	------	-------	-------------

- BRAM_PORTB: read interface to BRAM

signal	direction	type	width	description
--------	-----------	------	-------	-------------

4.1.2 Parameter

- PREVIOUS_LAYER_TYPE boolean: TRUE: conv2d, FALSE: dense
- PREVIOUS_LAYER_WIDTH integer: Row length of input matrix
- PREVIOUS_LAYER_HEIGHT integer: Column length of input matrix
- PREVIOUS_LAYER_CHANNEL integer: Row length of input matrix
- NEXT_LAYER_TYPE boolean: TRUE: conv2d, FALSE: dense
- NEXT_LAYER_WIDTH integer: Row length of input matrix
- NEXT_LAYER_HEIGHT integer: Column length of input matrix
- NEXT_LAYER_CHANNEL integer: Row length of input matrix

use extra parameter for dense or simply use width or height, discuss!

use extra parameter for dense or simply use width or height, discuss!

4.2 conv2d

Figure 3 shows the block diagram of a conv2d module. It uses k conv_channel modules to realise k output channels. All conv_channel modules get the same input vector X_{c_i} which consists of $n \cdot 3 \times 1$ vector, in which n is the number of input channels.

4.2.1 Interface

define input output interface

4.2.2 Parameter

- INPUT_CHANNEL_NUMBER : integer
- OUTPUT_CHANNEL_NUMBER : integer
- MATRIX_WIDTH: integer
- MATRIX_HEIGHT: integer

4.3 conv_channel

Figure 4 shows the block diagram of a conv_channel module. It uses n kernel_3x3 modules to realise n input channels. All kernel_3x3 modules get a different input vector $X_{c_{i1}}$ to $X_{c_{in}}$ which are 3×3 input matrices.

4.3.1 Interface

define input output interface

4.3.2 Parameter

- INPUT_CHANNEL_NUMBER : integer
- OUTPUT_CHANNEL_NUMBER : integer
- MATRIX_WIDTH: integer
- MATRIX_HEIGHT: integer

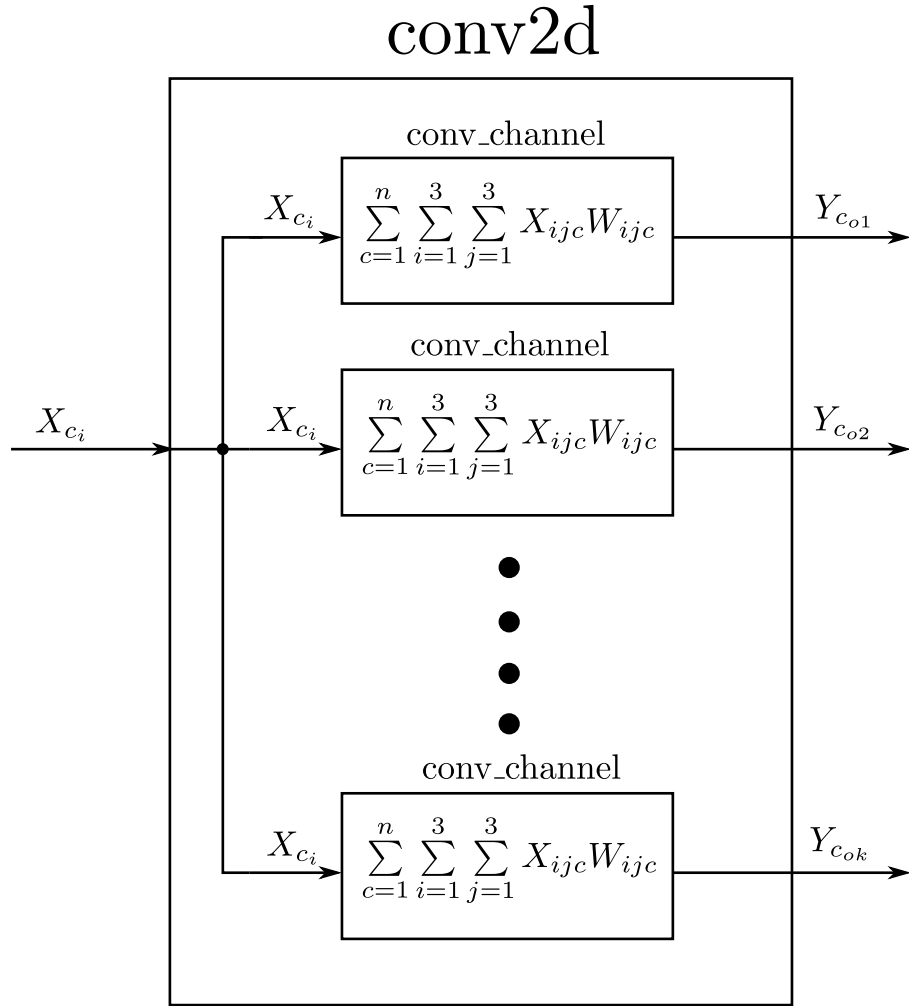
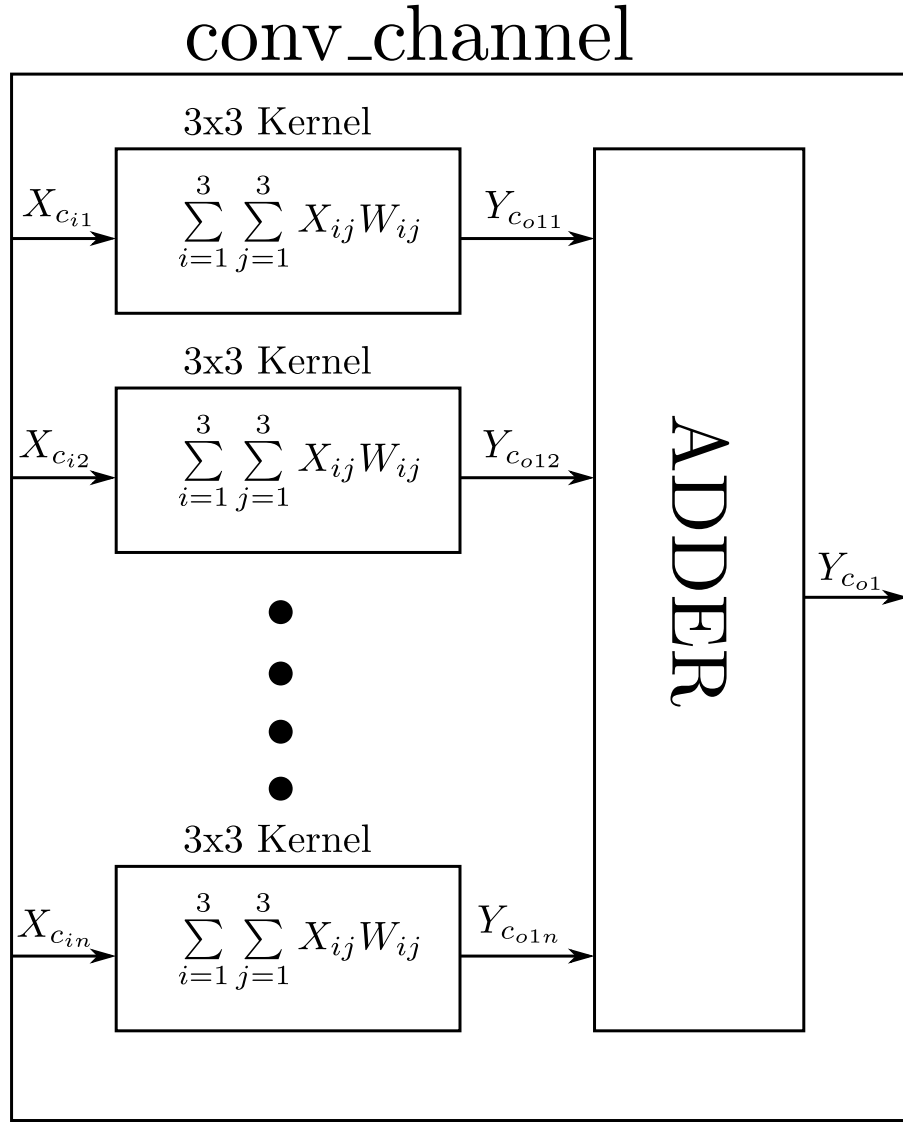


Figure 3: Conv2d block diagram. For each output channel a conv_channel module is used. k indicates the number of output channels.



Parameter:

- input channel number

Figure 4: conv_channel block diagram. For each input channel a kernel_3x3 module is used. n indicates the number of input channels.