

# High Level Synthesis of a trained CNN for handwritten digit recognition

Federico Serafini

`federico.serafini@studenti.unipr.it`

Embedded Systems

Università degli Studi di Parma

12/07/2022

# Outline

- 1 Introduction
- 2 SW implementation
- 3 High Level Synthesis
- 4 Results and Validation
- 5 Conclusions

# Basic concepts

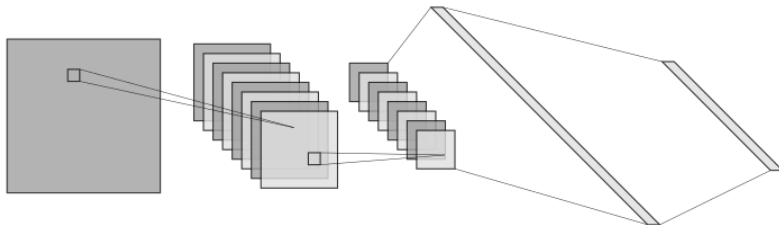
## Convolutional Neural Network (CNN)

Neural networks able to detect spatial structures (features) of the input through a special architecture based on:

- local receptive fields (convolution operation);
- shared weights (filters);
- local translation invariance (pooling operation).

⇒ Widely used in image-recongnition problems.

⇒ Highly-parallelizable problem.



# Basic concepts

## A single and ambitious objective

Overtake C performances through HW parallelism!

## Workflow

- ① Python:
  - ① model definition, training and evaluation;
  - ② export of network weights and architecture.
- ② C: replication of the network architecture.
- ③ Vitis HLS:
  - ① naive implementation (basic C synthesis);
  - ② stream and dataflow implementation.
- ④ Validation.

# Model definition and evaluation in Python

padding_layer_input	input:	[(None, 28, 28, 1)]
InputLayer		
float32	output:	[(None, 28, 28, 1)]

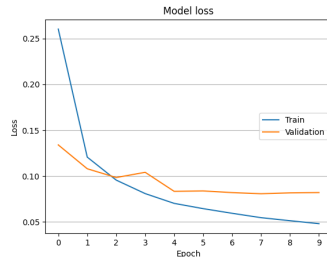
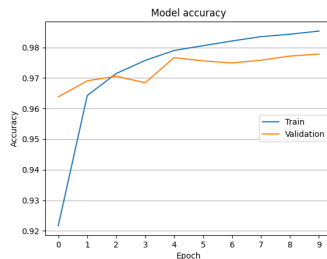
padding_layer	input:	(None, 28, 28, 1)
ZeroPadding2D		
float32	output:	(None, 30, 30, 1)

convolution_layer	input:	(None, 30, 30, 1)
Conv2D   relu		
float32	output:	(None, 28, 28, 8)

max_pooling_layer	input:	(None, 28, 28, 8)
MaxPooling2D		
float32	output:	(None, 14, 14, 8)

flatten_layer	input:	(None, 14, 14, 8)
Flatten		
float32	output:	(None, 1568)

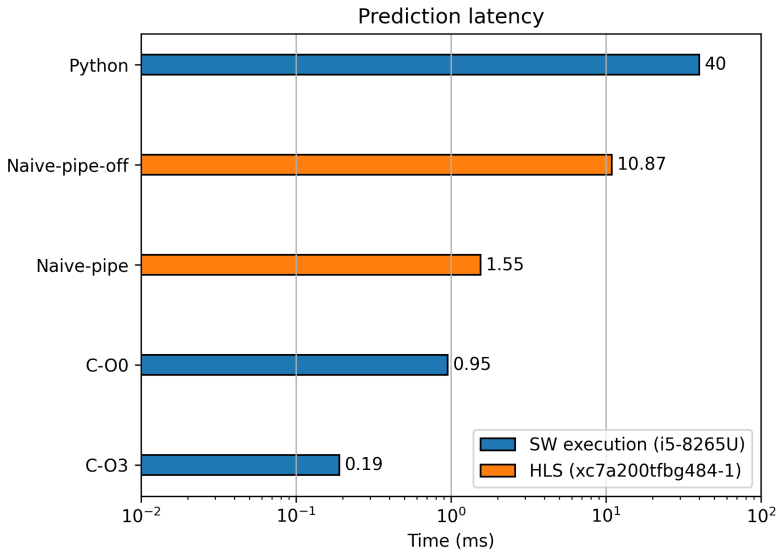
dense_layer	input:	(None, 1568)
Dense   softmax		
float32	output:	(None, 10)



# Network replication in C

```
1 void cnn
2 (
3     float img_in      [IMG_ROWS][IMG_COLS],
4     float prediction [DIGITS]
5 )
6 {
7     /***** Normalization and padding. *****/
8     float pad_img [PAD_IMG_ROWS][PAD_IMG_COLS] = { 0 };
9     normalization_and_padding(img_in, pad_img);
10
11     /***** Convolution layer. *****/
12     float features [FILTERS][IMG_ROWS][IMG_COLS] = { 0 };
13     // Convolution with relu as activation function.
14     convolutional_layer(pad_img, features);
15
16     /***** Max-pooling layer. *****/
17     float pool_features [FILTERS][POOL_IMG_ROWS][POOL_IMG_COLS] = { 0 };
18     max_pooling_layer(features, pool_features);
19
20     /***** Flatten layer. *****/
21     float flat_array [FLAT_SIZE] = { 0 };
22     flattening_layer(pool_features, flat_array);
23
24     /***** Dense layer. *****/
25     dense_layer(flat_array, prediction);
26 }
```

# Naive implementation: latency estimation



# Stream<sup>1</sup>

## Streaming data

A type of data transfer in which data samples are sent in sequential order starting from the first sample:

- a FIFO of infinite depth (*ap\_fifo*);
- no address management is required.

## Array implemented as FIFO interface

- Array must be only read or written, thus allowing a point-to-point connection.
- Program must follow first in, first out semantics (random access is not supported).
- If a stream is used to transfer data between tasks, consider a dataflow region where data streams from one task to the next.

---

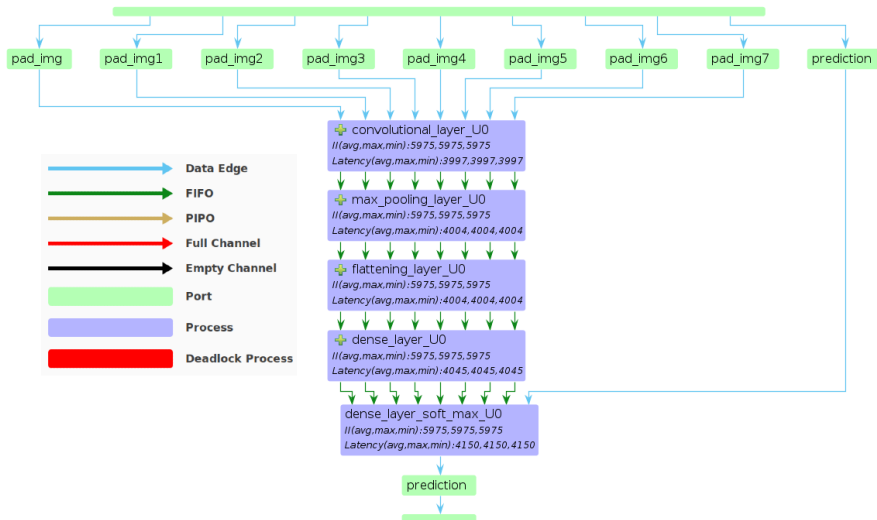
<sup>1</sup>Vitis-HLS User Guide



# Stream-dataflow implementation: code structure

```
1  #define FILTERS 8
2
3  void cnn
4  (
5      float img_in      [IMG_ROWS][IMG_COLS],
6      float prediction [DIGITS]
7  )
8  {
9      /***** Pre-processing the img_in. *****/
10
11      // Normalization and padding.
12
13
14      /***** Clone the normalized and padded image for FILTERS times. *****/
15
16      /*
17       *   Clone the normalized and padded image in order to
18       *   have an image for each parallel execution.
19       */
20
21
22      /***** Parallel executions start here. *****/
23
24      /*
25       *   Dataflow section with streams between tasks:
26       *   -convolution_layer;
27       *   -max_pooling_layer;
28       *   -flattening_layer;
29       *   -dense_layer.
30       */
31  }
```

# Dataflow view



# High level synthesis details report

## 'cnn' report

### Latency

#### Summary

Latency (cycles)		Latency (absolute)		Interval (cycles)		
min	max	min	max	min	max	Type
5821	5821	58.210 us	58.210 us	5822	5822	none

#### Detail

##### Instance

Instance	Module	Latency (cycles)		Latency (absolute)		Interval (cycles)		Type
		min	max	min	max	min	max	
grp_dataflow_section_fu_466	dataflow_section	3997	3997	39.970 us	39.970 us	3998	3998	dataflow

##### Loop

Loop Name	Latency (cycles)		Iteration Latency	Initiation Interval		Trip Count	Pipelined
	min	max		achieved	target		
-pad_for_rows_pad_for_cols	918	918	20	1	1	900	yes
-clone_for_rows_clone_for_cols	901	901	3	1	1	900	yes

## 'dataflow\_section' report

### Latency

#### Summary

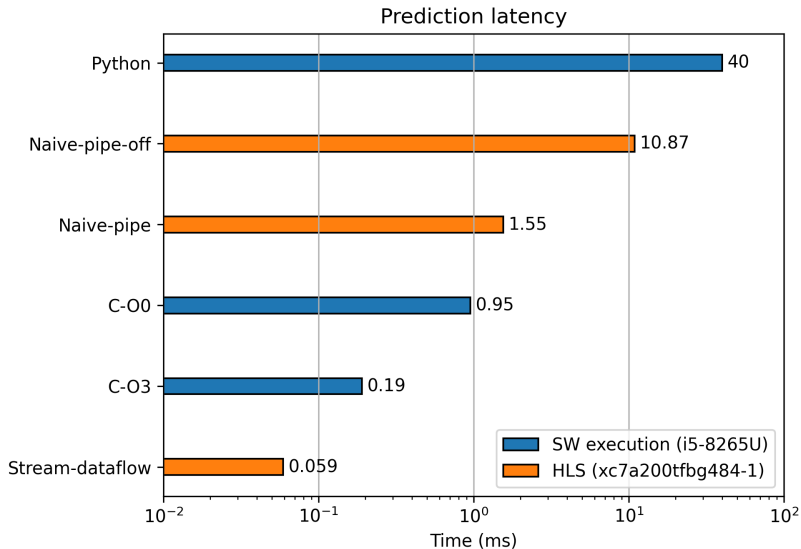
Latency (cycles)		Latency (absolute)		Interval (cycles)		
min	max	min	max	min	max	Type
3997	3997	39.970 us	39.970 us	3998	3998	dataflow

#### Detail

##### Instance

Instance	Module	Latency (cycles)		Latency (absolute)		Interval (cycles)		Type
		min	max	min	max	min	max	
convolutional_layer_U0	convolutional_layer	3997	3997	39.970 us	39.970 us	3998	3998	dataflow
dense_layer_U0	dense_layer	1998	1998	19.980 us	19.980 us	1999	1999	dataflow
dense_layer_soft_max_U0	dense_layer_soft_max	113	113	1.130 us	1.130 us	113	113	none
max_pooling_layer_U0	max_pooling_layer	793	793	7.930 us	7.930 us	794	794	dataflow
flattening_layer_U0	flattening_layer	198	198	1.980 us	1.980 us	199	199	dataflow

# Stream-dataflow implementation: latency estimation



# Validation

## Co-simulation

<b>Total predictions</b>	10
<b>Correct predictions</b>	100%

	<b>Min</b>	<b>Avg</b>	<b>Max</b>
<b>Latency (cycles)</b>	5975	5975	5975

## Export RTL with Vivado synthesis and place and route

	<b>BRAM</b>	<b>DSP</b>	<b>FF</b>	<b>LUT</b>
<b>Vitis HLS</b>	384	143	47201	37585
<b>Vivado</b>	224	143	38791	26753

	<b>Required</b>	<b>Post-synth</b>	<b>Post-impl</b>
<b>Clock period (ns)</b>	10	8.123	9.157

# Conclusions and future work

## Future work

- Smarter SW algorithm  $\implies$  faster HW accelerator.
- Fixed-point arithmetic  $\implies$  reduced area.
- Vitis HLS syntax constructs and libraries  $\implies$  smarter synthesis.

## Hardware or software implementation?

Further investigation is needed:

- ① apply improvements listed above;
- ② consider application domain requirements (for example the available HW and timing constraints);
- ③ consider also a mixed approach (HW and SW);
- ④ choose the cheapest solution satisfying the requirements.

Thank you for your attention.