# Demo CubeMX-AI: A CIFAR10 classification

Alessio Burrello,

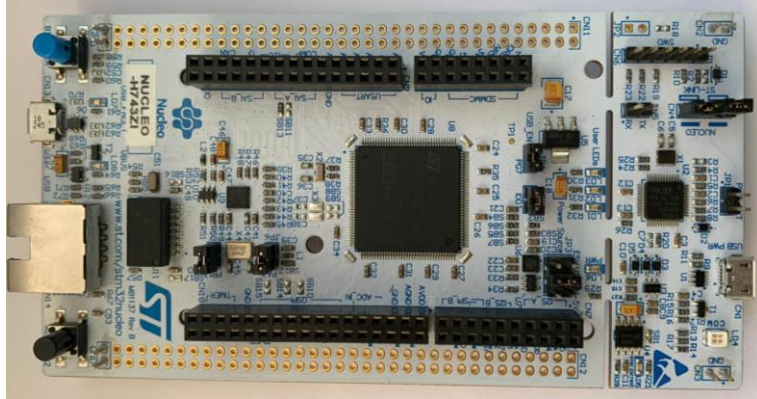alessio.burrello@unibo.it

Credits: Marcello Zanghieri,
Enrico Tabanelli

# Summary

- Creation of the python model and production of the saved model (.h5)

- Creation of the CubeMX project and validation on NUCLEO-H743ZI

- Creation of the application for network inference

# Material

- NUCLEO-H743ZI



- Cable to connect



- PC with Windows + VM with Linux
  - Python 3.x
  - STMCubeMX-AI
  - STM32CubeIDE

# Summary

- **Creation of the python network and production of the saved model (.h5)**

- Creation of the CubeMX project and validation on NUCLEO-H743ZI

- Creation of the application for network inference

# Python network generation: keras

## Training of the network

```
Epoch 78/100
50000/50000 [==============================] - 38s 765us/step - loss: 1.1609 - acc: 0.5928 - val_loss: 1.0284 - val_acc: 0.6407
Epoch 79/100
50000/50000 [==============================] - 37s 745us/step - loss: 1.1530 - acc: 0.5970 - val_loss: 1.0538 - val_acc: 0.6279
Epoch 80/100
50000/50000 [==============================] - 37s 731us/step - loss: 1.1509 - acc: 0.5976 - val_loss: 1.0668 - val_acc: 0.6264
Epoch 81/100
50000/50000 [==============================] - 38s 750us/step - loss: 1.1544 - acc: 0.5966 - val_loss: 1.0261 - val_acc: 0.6438
Epoch 82/100
50000/50000 [==============================] - 38s 759us/step - loss: 1.1492 - acc: 0.5994 - val_loss: 1.0408 - val_acc: 0.6353
Epoch 83/100
50000/50000 [==============================] - 39s 776us/step - loss: 1.1470 - acc: 0.6018 - val_loss: 1.0295 - val_acc: 0.6439
Epoch 84/100
50000/50000 [==============================] - 39s 780us/step - loss: 1.1445 - acc: 0.6021 - val_loss: 1.0320 - val_acc: 0.6437
Epoch 85/100
50000/50000 [==============================] - 38s 766us/step - loss: 1.1423 - acc: 0.6026 - val_loss: 1.0200 - val_acc: 0.6440
Epoch 86/100
50000/50000 [==============================] - 37s 736us/step - loss: 1.1420 - acc: 0.6030 - val_loss: 1.0344 - val_acc: 0.6391
Epoch 87/100
50000/50000 [==============================] - 36s 721us/step - loss: 1.1339 - acc: 0.6049 - val_loss: 1.0355 - val_acc: 0.6385
Epoch 88/100
50000/50000 [==============================] - 36s 726us/step - loss: 1.1280 - acc: 0.6076 - val_loss: 1.0045 - val_acc: 0.6500
Epoch 89/100
50000/50000 [==============================] - 36s 720us/step - loss: 1.1349 - acc: 0.6072 - val_loss: 1.0171 - val_acc: 0.6443
Epoch 90/100
50000/50000 [==============================] - 32s 635us/step - loss: 1.1269 - acc: 0.6088 - val_loss: 1.0156 - val_acc: 0.6450
Epoch 91/100
50000/50000 [==============================] - 32s 644us/step - loss: 1.1258 - acc: 0.6085 - val_loss: 1.0095 - val_acc: 0.6429
Epoch 92/100
50000/50000 [==============================] - 34s 675us/step - loss: 1.1222 - acc: 0.6118 - val_loss: 1.0126 - val_acc: 0.6464
Epoch 93/100
50000/50000 [==============================] - 35s 697us/step - loss: 1.1242 - acc: 0.6126 - val_loss: 1.0471 - val_acc: 0.6336
Epoch 94/100
50000/50000 [==============================] - 41s 810us/step - loss: 1.1157 - acc: 0.6146 - val_loss: 1.0123 - val_acc: 0.6432
Epoch 95/100
50000/50000 [==============================] - 38s 767us/step - loss: 1.1165 - acc: 0.6139 - val_loss: 1.0133 - val_acc: 0.6482
Epoch 96/100
50000/50000 [==============================] - 36s 726us/step - loss: 1.1177 - acc: 0.6155 - val_loss: 1.0527 - val_acc: 0.6300
Epoch 97/100
50000/50000 [==============================] - 36s 729us/step - loss: 1.1105 - acc: 0.6163 - val_loss: 0.9951 - val_acc: 0.6513
Epoch 98/100
50000/50000 [==============================] - 39s 775us/step - loss: 1.1149 - acc: 0.6135 - val_loss: 1.0062 - val_acc: 0.6481
Epoch 99/100
50000/50000 [==============================] - 35s 699us/step - loss: 1.1095 - acc: 0.6142 - val_loss: 1.0103 - val_acc: 0.6473
Epoch 100/100
50000/50000 [==============================] - 34s 675us/step - loss: 1.1100 - acc: 0.6166 - val_loss: 1.0039 - val_acc: 0.6520
```

## How to train

- Launch CIFAR10 network example, that you find in the folder
- Change the following parameters, to generate more models:

*batch_size = 32*
*num_classes = 10*
*epochs = 200*
*data_augmentation = True*
*num_predictions = 20*
*save_dir = os.path.join(os.getcwd(), 'saved_models')*
*model_name = 'keras_cifar10_trained_model_200_epochs.h5'*

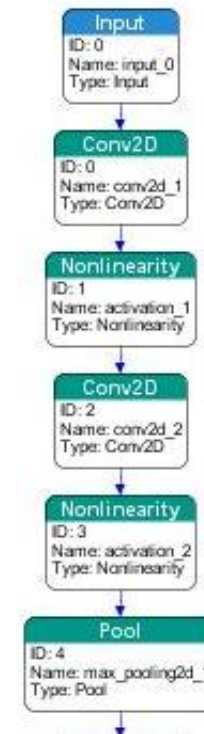# Python network generation: keras

## Generated model:

- .h5 file, with topology + weights
- a CNN network with conv2d, poolings, Relu, Dropout and Dense Layers

Full network contains:
- 3 blocks identical to the presented one;
- 2 final Fully connected layers

Block 1 as visualized from CubeMX-AI

# Summary

- Creation of the python network and production of the saved model (.h5)

- **Creation of the CubeMX project and validation on NUCLEO-H743ZI**

- Creation of the application for network inference

# STM32CubeMX: New project



- Create a new project
- Select **Board Selector**
- Select **NUCLEO-H743ZI**
*Select a different board if you have a different one.*

# STM32CubeMX: add AI



- Use the *Additional software* to add a validation application of a Neural network.
- Validate the execution of the NN on the NUCLEO-H7.
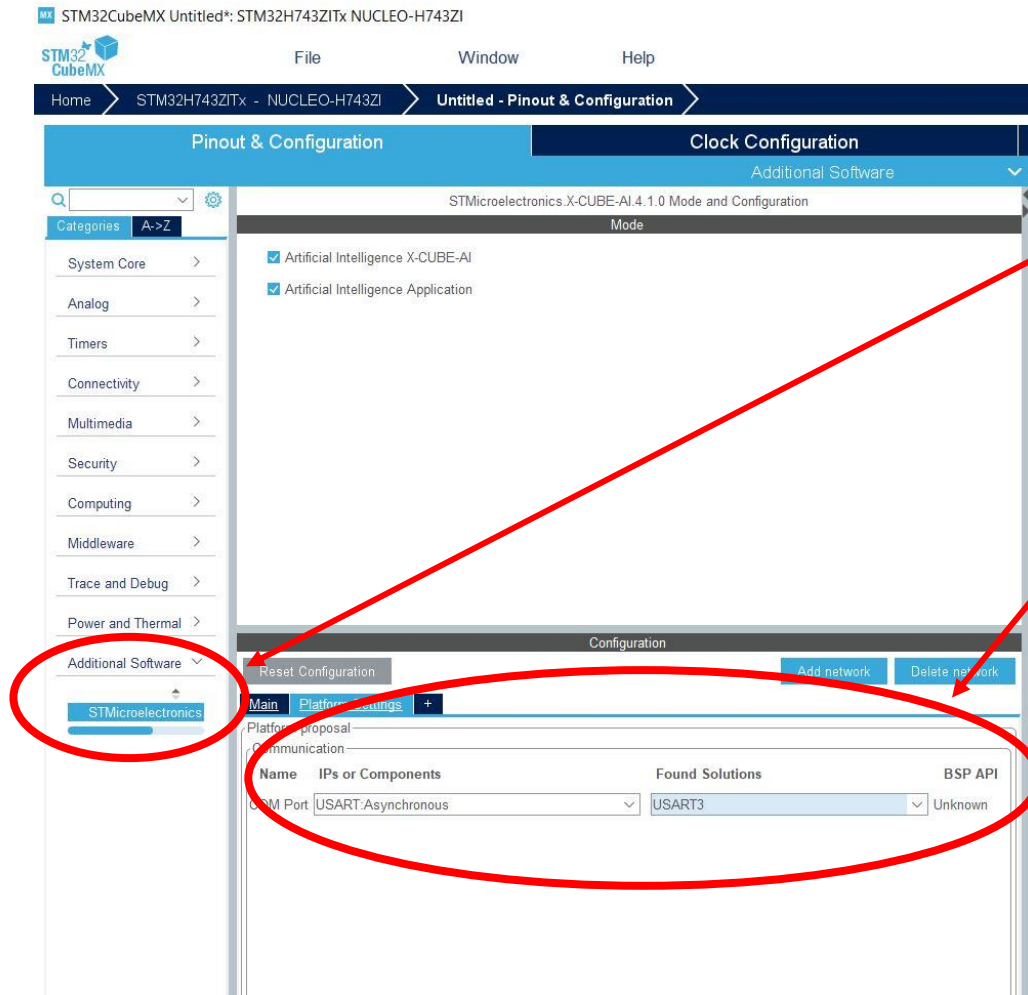
# STM32CubeMX: add AI



- Select *Validation* to validate your application (check memory and accuracy of the result → MSE compared to PC model)

# STM32CubeMX-AI: Configuration of the network



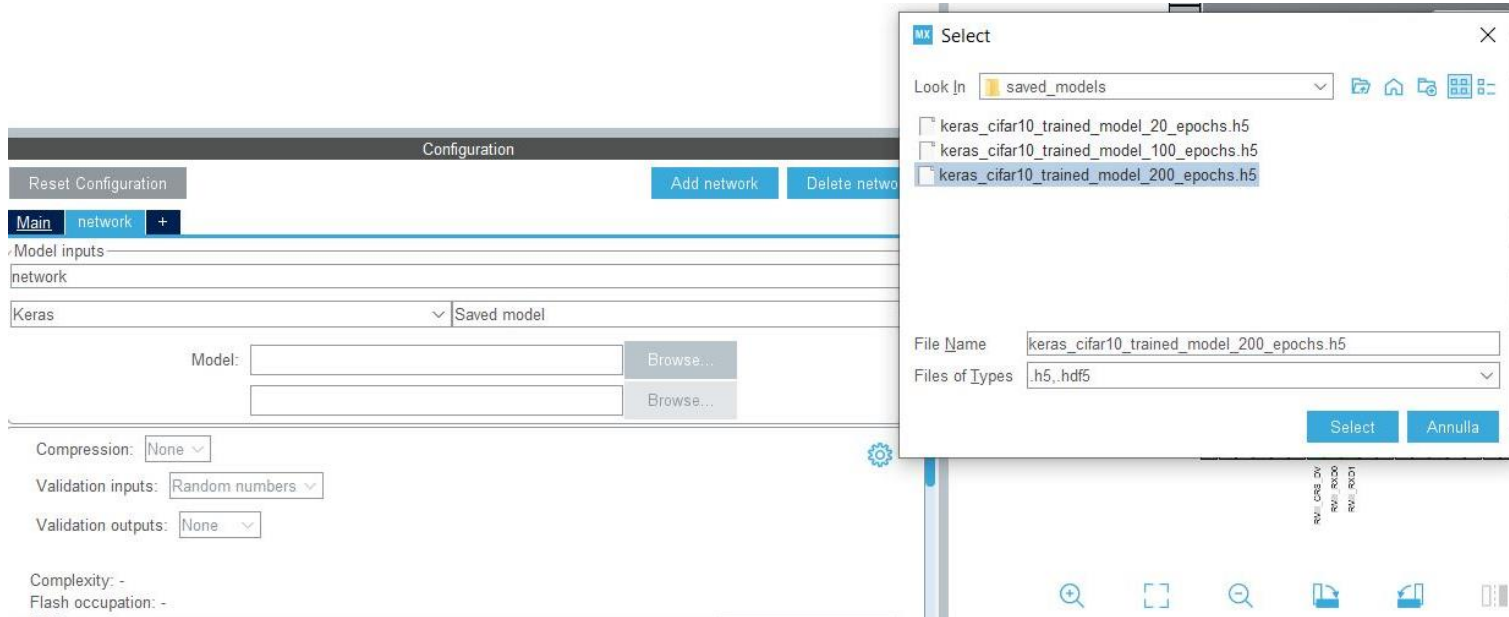- Additional software configured

# STM32CubeMX-AI: Configuration of the network



- Additional software configured

- Configure the USART3 for validation in Platform settings.
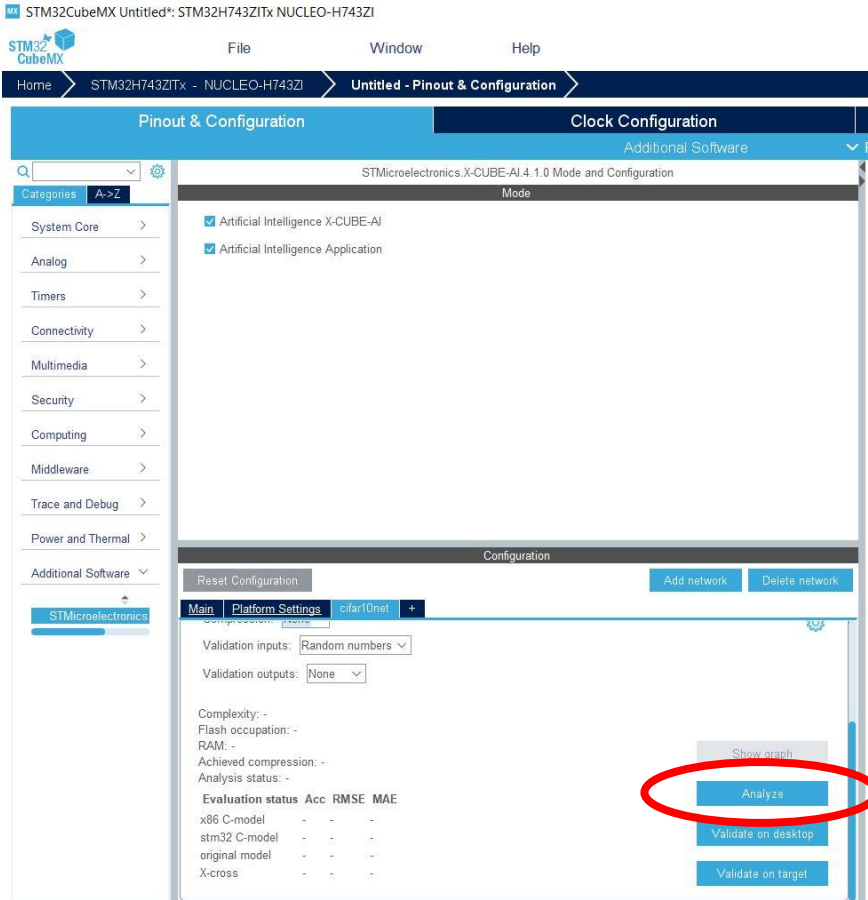  - If not automatically activated, remember to activate USART3 in *Connectivity*

**N.B.** *USART3 is the one connected to the STLINK for NUCLEO-H7. Check the right one on your board.*
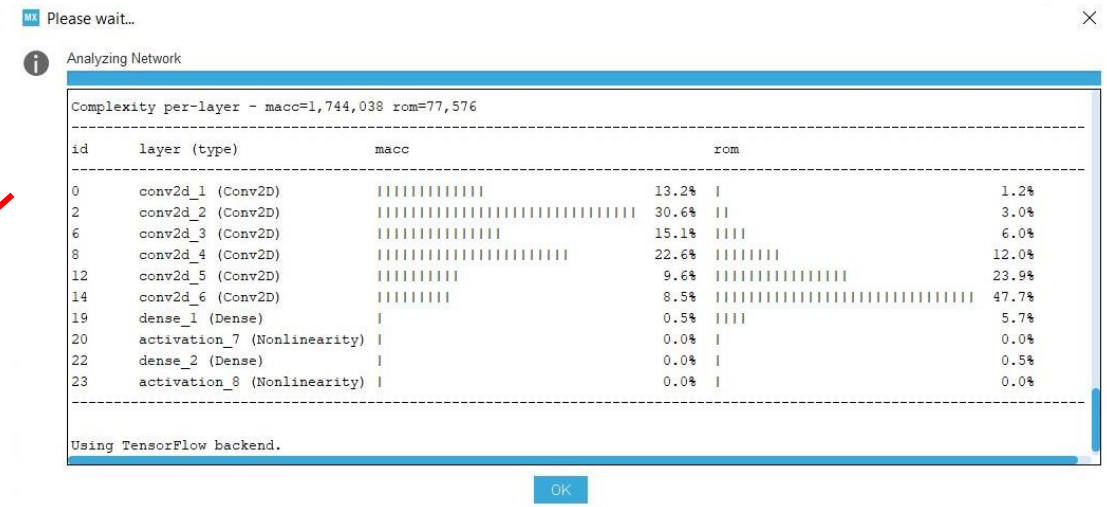
# STM32CubeMX-AI: Add a network



- Select *keras, saved model*.
- Change name of Model input to ***cifar10net***
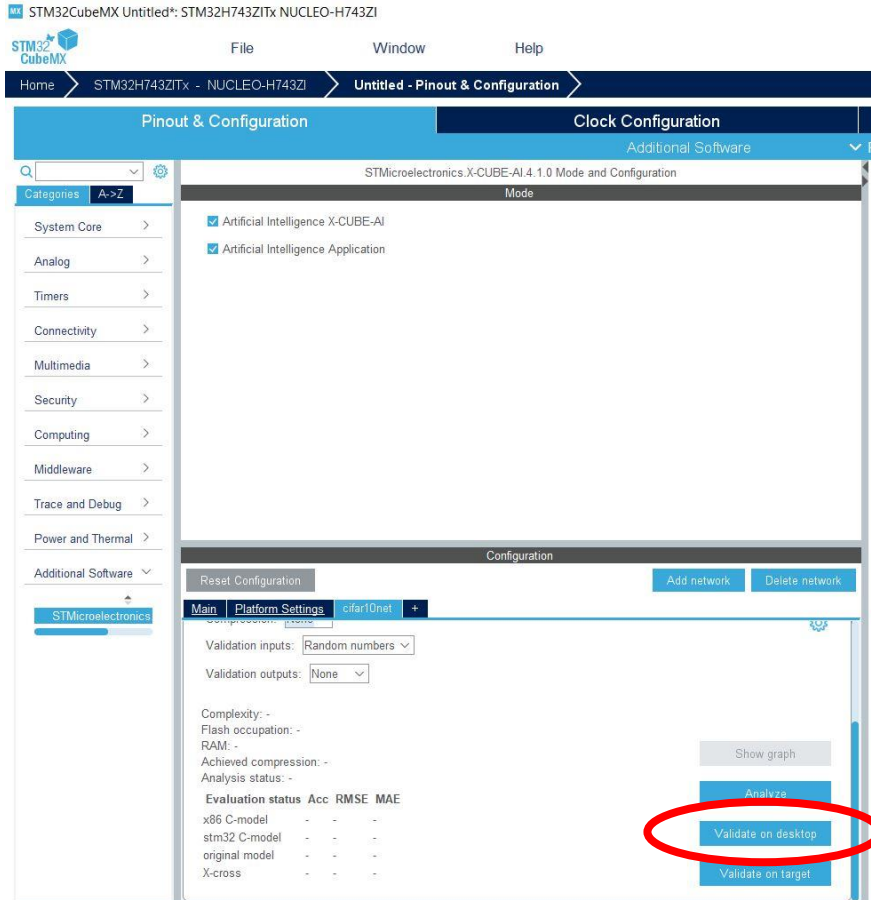- Import a network from saved models in github (or a generic .h5 network)
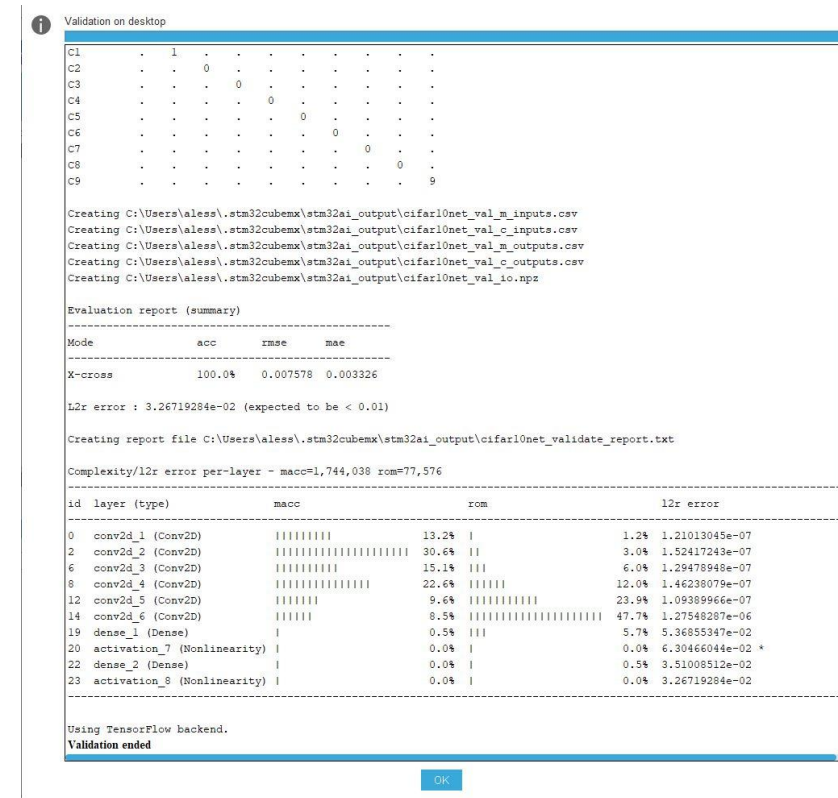
# STM32CubeMX-AI: actions



**Analyze**. Check the correctness of the network. You can choose a level of compression (None, 4, 8)
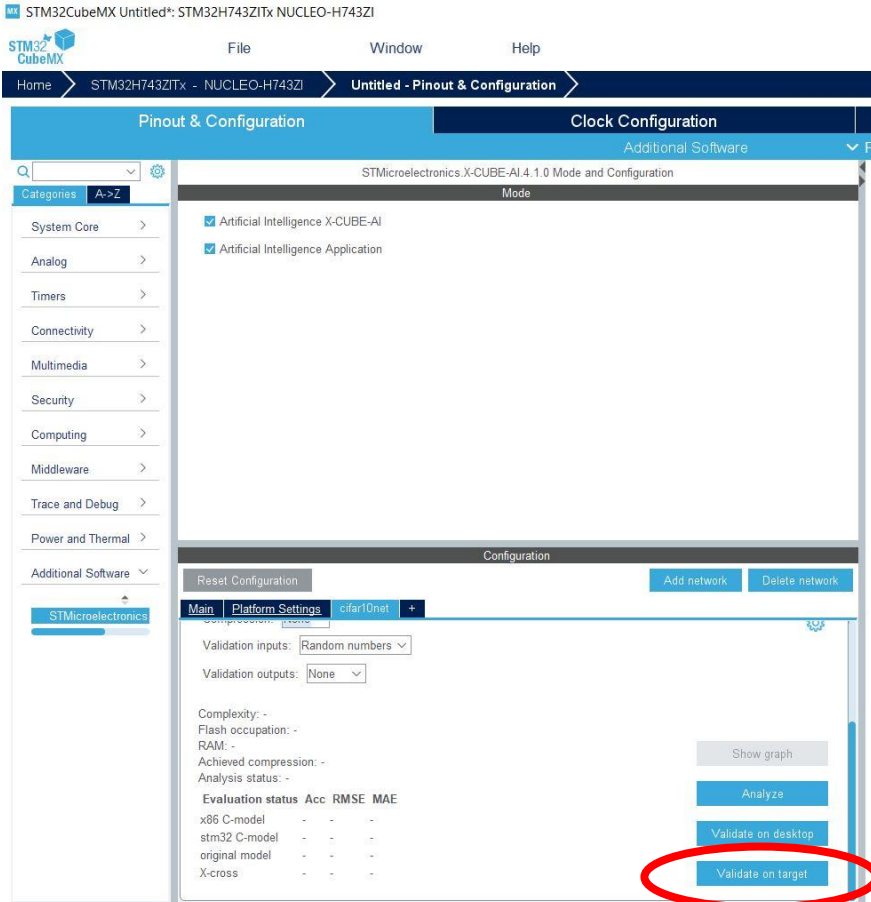
# STM32CubeMX-AI: actions
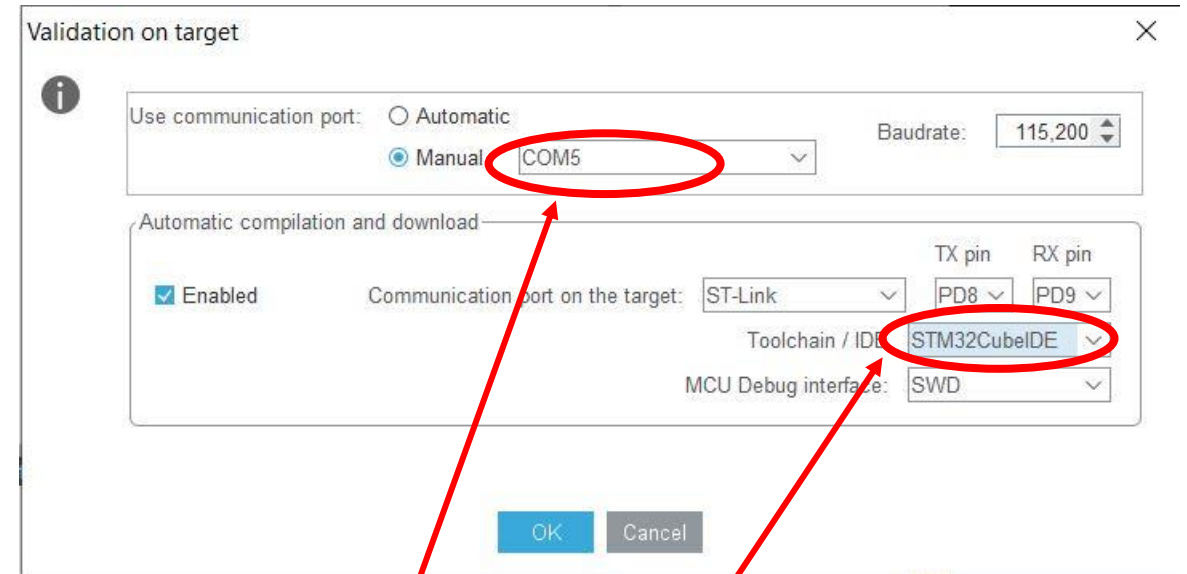


## Validate on desktop

# STM32CubeMX-AI: actions



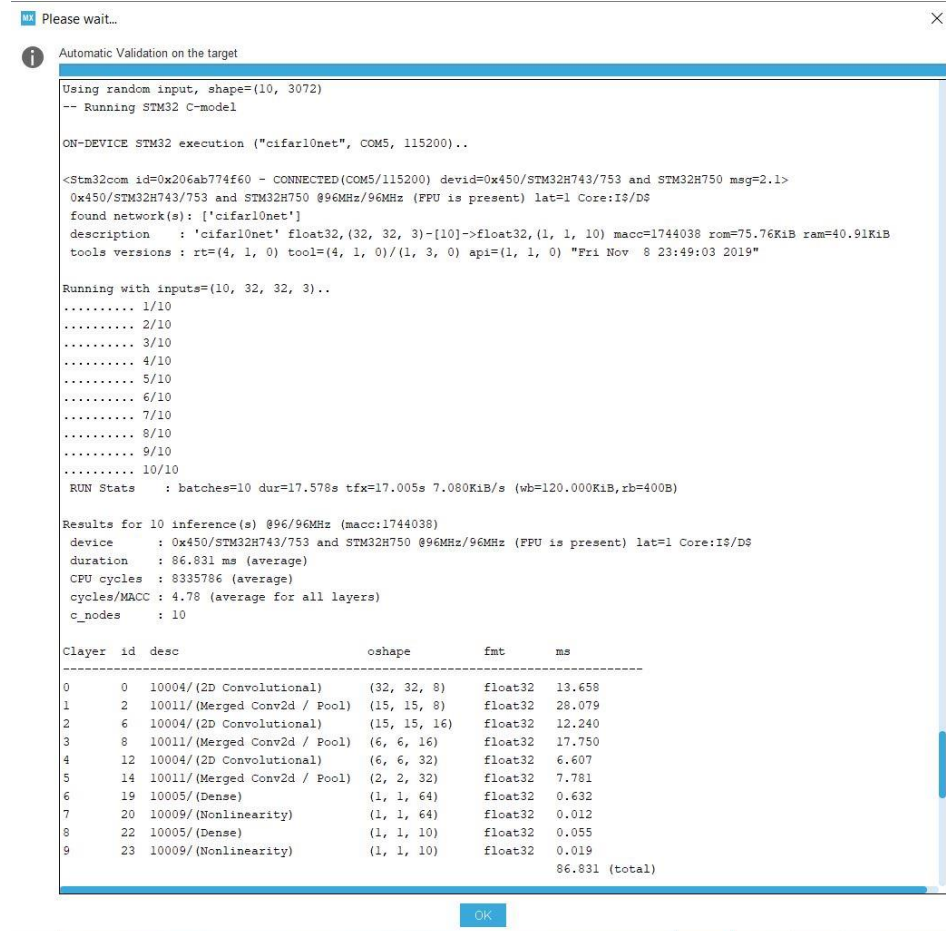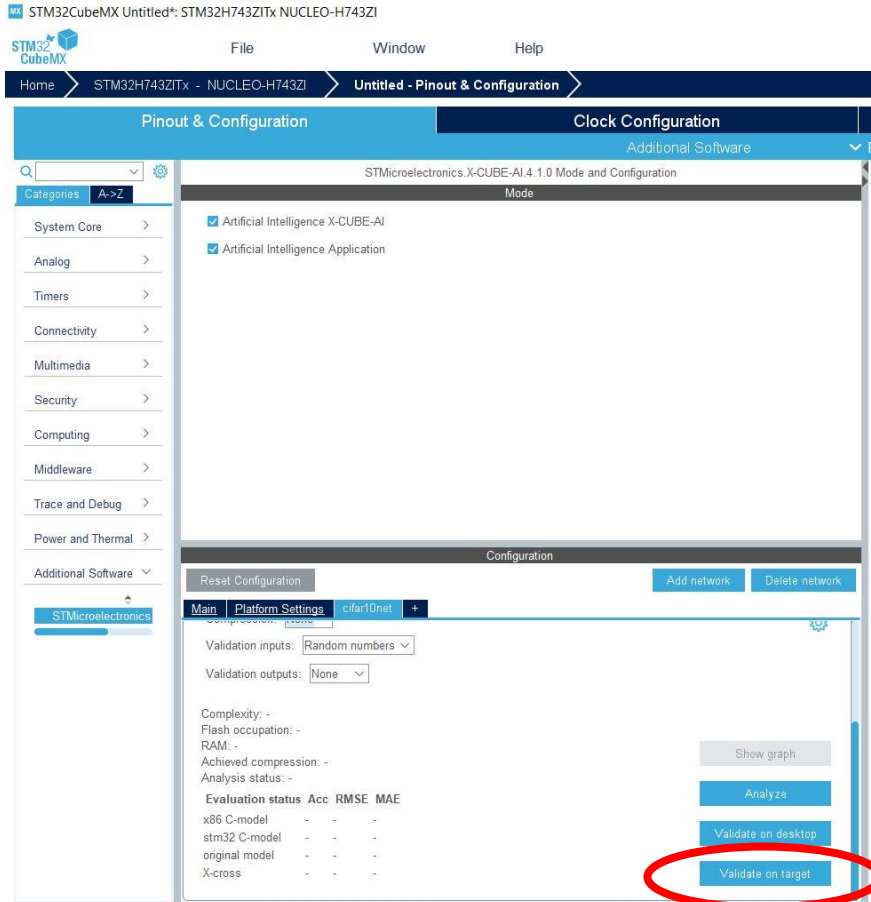**Validate on target:** need to have the NUCLEO-H7 plugged.



*Select correct IDE and COM*

# STM32CubeMX-AI: actions



**Validate on target:** need to have the NUCLEO-H7 plugged.
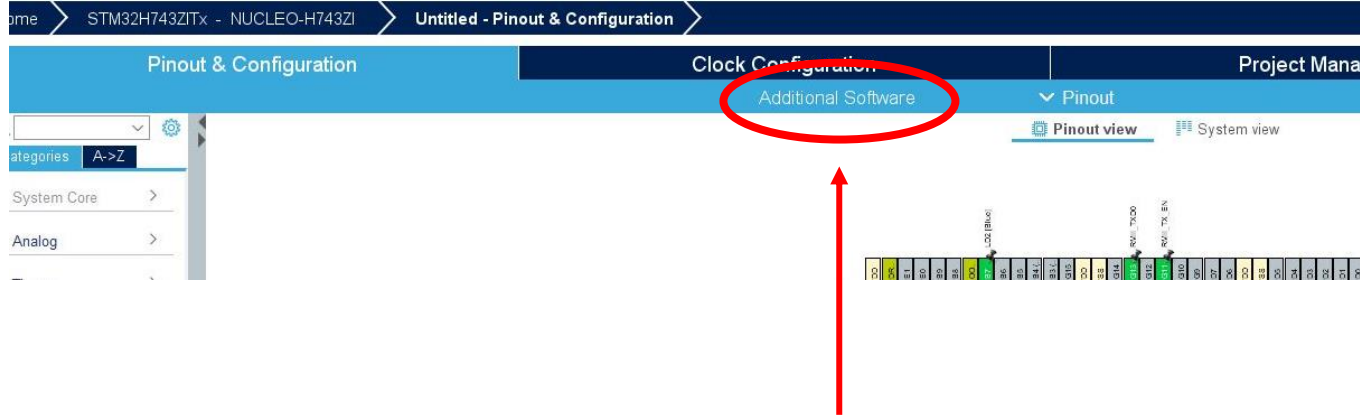
# Summary

- Creation of the python network and production of the saved model (.h5)

- Creation of the CubeMX project and validation on NUCLEO-H743ZI

- **Creation of the application for network inference**

# STM32CubeMX-AI: application generation



- Change *Validation* to *Application template*. The application template allows to Generate the code using the *Generate code* button.

# STM32CubeMX-AI: application generation



- Go to *Project Manager*, give a name to the project (three projects are available on Dropbox)

# STM32CubeMX-AI: application generation



- change the IDE to **STM32CubeIDE**

# STM32CubeMX-AI: application generation



- *Generate Code*

# STM32CubeMX-AI Application: add Custom code

- Add Custom code to communicate with the NUCLEO Board and classify with the network
    - USART Receive to receive the image from the PC
    - USART Transmit to send the results back to the PC


1. File to be modified → app_x-cube-ai.c
2. Modifications:
    - Usart declaration (we used USART3 since it is the one linked to the STLINK)
    - Usart initialization
    - Usart communication

# STM32CubeMX-AI Application: add Custom code

```c
/* USER CODE BEGIN includes */
#include <stdio.h>
UART_HandleTypeDef huart3;

/* USER CODE END includes */
```

1. Usart declaration

# STM32CubeMX-AI Application: add Custom code

```
/* USER CODE BEGIN includes */
#include <stdio.h>
UART_HandleTypeDef huart3;

/* USER CODE END includes */
```

1. Usart declaration
2. Usart init function declaration (copy from *main.c* file)

```
static void MX_USART3_UART_Init(void)
{

    /* USER CODE BEGIN USART3_Init 0 */

    /* USER CODE END USART3_Init 0 */

    /* USER CODE BEGIN USART3_Init 1 */

    /* USER CODE END USART3_Init 1 */
    huart3.Instance = USART3;
    huart3.Init.BaudRate = 115200;
    huart3.Init.WordLength = UART_WORDLENGTH_8B;
    huart3.Init.StopBits = UART_STOPBITS_1;
    huart3.Init.Parity = UART_PARITY_NONE;
    huart3.Init.Mode = UART_MODE_TX_RX;
    huart3.Init.HwFlowCtl = UART_HWCONTROL_NONE;
    huart3.Init.OverSampling = UART_OVERSAMPLING_16;
    huart3.Init.OneBitSampling = UART_ONE_BIT_SAMPLE_DISABLE;
    huart3.Init.ClockPrescaler = UART_PRESCALER_DIV1;
    huart3.AdvancedInit.AdvFeatureInit = UART_ADVFEATURE_NO_INIT;
```

# STM32CubeMX-AI Application: add Custom code

```
/* USER CODE BEGIN includes */
#include <stdio.h>
UART_HandleTypeDef huart3;


/* USER CODE END includes */
```

```
static void MX_USART3_UART_Init(void)
{

    /* USER CODE BEGIN USART3_Init 0 */

    /* USER CODE END USART3_Init 0 */

    /* USER CODE BEGIN USART3_Init 1 */

    /* USER CODE END USART3_Init 1 */
    huart3.Instance = USART3;
    huart3.Init.BaudRate = 115200;
    huart3.Init.WordLength = UART_WORDLENGTH_8B;
    huart3.Init.StopBits = UART_STOPBITS_1;
    huart3.Init.Parity = UART_PARITY_NONE;
    huart3.Init.Mode = UART_MODE_TX_RX;
    huart3.Init.HwFlowCtl = UART_HWCONTROL_NONE;
    huart3.Init.OverSampling = UART_OVERSAMPLING_16;
    huart3.Init.OneBitSampling = UART_ONE_BIT_SAMPLE_DISABLE;
    huart3.Init.ClockPrescaler = UART_PRESCALER_DIV1;
    huart3.AdvancedInit.AdvFeatureInit = UART_ADVFEATURE_NO_INIT;
```

1. Usart declaration
2. Usart init function declaration (copy from *main.c* file)
3. Usart initialization in MX_X_CUBE_AI_Init

```
     */
 void MX_X_CUBE_AI_Init(void)
 {
     /* USER CODE BEGIN 0 */
     /* Activation/working buffer is allocated as a static memory chunk
      * (bss section) */
     AI_ALIGNED(4)
     static ai_u8 activations[AI_CIFAR10NET_DATA_ACTIVATIONS_SIZE];
     MX_USART3_UART_Init();
     aiInit(activations);
     /* USER CODE END 0 */
 }
```

# STM32CubeMX-AI Application: add Custom code

1. Usart declaration
2. Usart init function declaration (copy from *main.c* file)
3. Usart initialization in MX_X_CUBE_AI_Init
4. MX_X_CUBE_AI_Process → change the input data stream from random data to USART transmitted data

```c
/* ------------------------------------------- */
/* Data generation and Pre-Process             */
/* ------------------------------------------- */
/* ------------------------------------------- */
ai_i8 in_data_temp[32*32*3];
while(HAL_UART_Receive(&huart3, in_data_temp, 32*32*3, 0xFFFF)!=HAL_OK);
for (ai_size i=0;  i < AI_CIFAR10NET_IN_1_SIZE; i++ )
{
    ai_float val =  (ai_float)in_data_temp[i];
    if (val<0)
        val = val+256.0;
    ((ai_float *)in_data)[i] = val/255.0;
}
/* Perform the inference */
aiRun(in_data, out_data);
res = 0;
ai_float pred = ((ai_float *)out_data)[0];
for (uint8_t i = 1; i<10;i++)
{
    if ((((ai_float *)out_data)[i] > pred)
    {
        pred = ((ai_float *)out_data)[i];
        res = i;
    }
}
while(HAL_UART_Transmit(&huart3, &res, sizeof(uint8_t), 0xFFFF)!=HAL_OK);
char msg[5] = "\n";
while(HAL_UART_Transmit(&huart3, msg, strlen(msg), 0xFFFF)!=HAL_OK);
        /* Post-Process - process the output buffer */
```

# STM32CubeMX-AI Application: load the code

- Load the code using the STM32CubeIDE on the board, which will be ready to run the network, waiting for an input image.