



AI:How to use transfer learning to perform image classification on STM32

AI:How to use transfer learning to perform image classification on STM32



Contents

1. AI:How to use transfer learning to perform image classification on STM32	3
2. Main Page	3



Stable: 25.05.2021 - 08:17 / Revision: 25.05.2021 - 08:17

The content format pdf is not supported by the content model wikitext.

[Return to Main Page](#)

Stable: 25.05.2021 - 08:18 / Revision: 25.05.2021 - 08:17

You do not have permission to edit this page, for the following reasons:

- The action you have requested is limited to users in one of the groups: **Administrators**, **Editors**, **Reviewers**, Selected_editors, ST_editors.
- The action "Read pages" for the draft version of this page is only available for the groups ST_editors, ST_readers, sysop, Selected_editors, Selected_readers, reviewer

You can view and copy the source of this page.



This article presents a video on how to use a technique called "Transfer learning" to quickly train a deep learning model in order to classify images. This video teaches you how to use ST ecosystem to build a computer vision application from the ground up. It describes how the FP-AI-VISION1 allows easy image collection with an STM32 Discovery kit, how to use transfer learning with Tensorflow to quickly train an image classification model, and eventually how to use STM32Cube.AI to convert this model into optimized code for ST microcontrollers. The Jupyter notebook used in this video is [https://github.com/STMicroelectronics/stm32ai/blob/master/AI_resources/VISION/transfer_learning/TransferLearning.ipynb available on Github]. It can be opened with Colab.   <https://www.youtube.com/watch?v=NDshmSH7WnA&list=PLnMKNibPkDnG9IC5NI9vJg1CKMAO1kODW> How to use transfer learning to perform image classification on STM32]]  The article below describes step by step the video tutorial. For another complete example on how to use your own models with the FP-AI-VISION1 and Teachable Machine [How to use Teachable Machine to create an image classification application on STM32|check out this article]]. In addition, [<https://www.youtube.com/watch?v=ftM2D6qgrtl> this series of videos] gives a good overview of the Vision function package. == Dataset generation == First you need a dataset of images representative of your classes. It can be generated either by following the steps described below or by using an already existing dataset (see related [https://github.com/STMicroelectronics/stm32ai/tree/master/AI_resources/VISION/transfer_learning Github page] for pasta datasets to jump start, note that it has been recorded with a fix distance of 13 cm and good light conditions). It is a good idea to use the same sensor for the dataset collection as the one that will be used for inference. In the present example, the [<https://www.st.com/en/development-tools/b-cams-omv.html> B-CAMS-OMV] board is used as camera module. Download the FP-AI-VISION1 function package from [<https://www.st.com/en/development-tools/b-cams-omv.html> <https://www.st.com/en/embedded-software/fp-ai-vision1.html>]. Using the USB Webcam application provided within FP-AI-VISION1, the board and the camera can be used as a webcam. First plug your board onto your computer using the "ST-LINK" port. Make sure the "JP6" jumper is set to "ST-LINK".   Discovery board ST-LINK setup]]  After plugging the USB cable onto your computer, the board appears as a mounted device. The binary is located under `FP-AI-VISION1_V3.0.0\Projects\STM32H7471-DISCO\Applications\USB_Webcam\Binary`. Drag and drop it onto the board mounted device. This flashes the binary on the board. Unplug the board, change the "JP6" jumper to the "HS" position, and plug your board using the "USB OTG" port.   Discovery board USB OTG setup]]  Launch the Windows camera application (it can be found by searching for "camera" in the Windows search bar). If your computer has already a camera, change the source of the image feed by clicking the "flip camera" button. The capture can be done either in "VGA" or "QVGA", the resolution can be changed with the menu `settings -> Photos -> Quality`. The images will be resized to the input resolution of the neural network in the Python notebook. Take enough picture samples to fill your dataset. It must be structured as follows:

```
Dataset cats
cat0001.jpg cat0002.jpg ... dogs horses
```

 The dataset must be large enough. The one used for the video includes 1408 files belonging to 8 classes. It is recommended to use at least 100 pictures per class. == Model generation == The Python notebook used to generate the model can be found at https://github.com/STMicroelectronics/stm32ai/blob/master/AI_resources/VISION/transfer_learning/TransferLearning.ipynb It can be open either with Google Colab (as in the present example) or by using any Python environment. Just click (or `Ctrl + Click`) the `Open in Colab` icon at the top of the notebook on Github. This takes you to the Colab notebook project. Zip your dataset folder and import it to the Colab project.   Colab import]]  You can run the script either step by step or run the whole script by pressing `Ctrl + F9`. All the steps are described in details in the script but what they

do is straightforward :

- * Import all the necessary libraries.
- * Unzip the dataset.
- * Import the dataset (this step also automatically resizes images to match the model input size).
- * Display samples from the dataset
- * Perform data augmentation by adding duplicate images with random contrast/translation/rotation/zoom.
- * Display new samples.
- * Normalize data.
- * Import a pretrained `MobileNet` model without the last layer and add new ones at the end.
- * Train the model.
- * Plot training loss history.
- * Quantize the model to reduce its size.
- * Export the quantized model. The trained and quantized model appears as `model_quant.tflite`.

== Model deployment == This section describes how to generate a model that can be used on STM32 platforms, starting from a model produced by a high-level library such a Keras/Tensorflow. It also describes how to use your model in an already existing project. === Model generation for STM32 === {{Warning|Warning, the latest version of FP-AI-VISION1 is based on STM32Cube.AI 6.0.0. Make sure your are using this version.}} First convert the `tflite` model to be used on a STM32 platform. There are two ways to achieve this goal, one by using STM32CubeMX and generating an entire project (recommended at first), the other by using the command line interface to generate only the required files. ==== Project generation with STM32CubeMX (Option 1) ==== Start STM32CubeMX and create a new project. In this example, the STM32H7471-DISCO board is used. Select



Start STM32CubeMX and create a new project. In this example, the STM32H747-DISCO board is used. Select the board in the `Board Selector` and start the project. A new window opens showing the board configuration. Click `Software Packs -> Select Components`.

[[File:CubeMXScreenShot.png|center|Adding components to STM32CubeMX]]

 Select the components for the Arm^{®} Cortex^{®}-M7 core under STMicroelectronics X-CUBE-AI, check the core box, select `Application` as application, then click Ok.

[[File:CubeMXScreenShot2.png|center|Components selector]]

 The stm32-ai library will be download if it has not been downloaded yet. When it is complete, click Ok. Go to `Software Packs (on the left) -> STMicroelectronics XCUBE AI`. A pop up window opens to configure peripherals and clock. Click Yes. Click `Add network`, call your network `network`, select TFLite from the scrollable menu, and browse to the last section to find the generated model. Once the model is selected, click `Analyze`. This provides a quick overview of the model memory footprint. Click the `Show graph` button for more in-depth information on the model topology and RAM memory usage layer by layer including activation buffers.

[[File:Graph1.png|center|Memory Usage windows]]

 To customize how the model is stored in memory, click the gear icon (`Advanced Settings`) and specify if you want to use external Flash memory and/or RAM. Check the `Use activation buffer for input buffer` to make them share the same buffer in order to reduce memory usage.

[[File:CubeMXScreenShot3.png|950px|center|X-Cube-AI model loading]]

 Click `Analyze` again. You can see that the RAM occupation has been significantly reduced (refer to the memory graph for more details).

[[File:Graph2.png|center|Memory Usage window after enabling input buffer space sharing]]

 Check that the model works properly by validating it either on the PC or on the target. To generate the code for STM32CubeIDE (or Arm^{®} KEIL^{®} or IAR SYSTEMS^{®}) go to `Project manager`, specify a project name and location. Select STM32CubeIDE as toolchain, then click `GENERATE CODE` on the top right corner.

[[File:CubeMXScreenShot4.png|950px|center|Project manager window]]

 This generates the necessary code for the whole project. In the present example, only the files related to the neural network are required. They are located under `YourProjectName\CM7\X-CUBE-AI\App` and are called `network.c`, `network.h`, `network_data.c`, `network_data.h` and `network_config.h`. They represent the model and its weight. ===== Command line interface (Option 2) ===== The other way to generate the necessary files is to use the command line interface. The main advantage is to produce only the required files. The documentation related to X CUBE AI is located under `C:\Users\<USERNAME>\STM32Cube\Repository\Packs\STMicroelectronics\X-CUBE-AI\6.0.0\Documentation\`. The CLI is called `stm32ai.exe` for Windows. It is located under `C:\Users\<USERNAME>\STM32Cube\Repository\Packs\STMicroelectronics\X-CUBE-AI\6.0.0\Utilities\windows` and `$HOME/STM32Cube/Repository/Packs/STMicroelectronics/X-CUBE-AI/6.0.0/Utilities\` for UNIX^{®}-based systems. If there is no `utilities` folder, create a CUBE MX project and add Core AI in the Software Package (see section Project generation). This downloads all the necessary files. For ease of usage, add the folder to your Windows path:

```
{{PC$}} set CUBE_FW_DIR=C:\Users\&lt;USERNAME&gt;\STM32Cube\Repository {{PC$}} set X_CUBE_AI_DIR=%CUBE_FW_DIR%\Packs\STMicroelectronics\X-CUBE-AI\6.0.0 {{PC$}} set PATH=%X_CUBE_AI_DIR%\Utilities\windows;%PATH%
```

 Once added, go to the folder where you wish to generate the network files in command line (using PowerShell or similar) and run the following command:

```
{{PC$}} stm32ai.exe generate -m .\model_quant.tflite --allocate-inputs # (change the path to the model if needed).
```

 This creates all the necessary files under the `stm32ai_output` folder. === Model integration === Now that the C files are generated for your model, integrate them in an already existing project. In this example, the application is based on the "Person detection" AI function package project. Launch the STM32CubeIDE project with the `.project` located under `FP-AI-VISION1_V3.0.0\FP-AI-VISION1_V3.0.0\Projects\STM32H747-DISCO\Applications\PersonDetection\MobileNetv2_Model\STM32CubeIDE`. Double click it to open STM32CubeIDE. Select the `CM7` subproject and build it. The resulting workspace is shown below:

[[File:CubeIDEScreenShot.png|center|STM32CubeIDE Project window]]

 Replace the model files (`network.c`, `network.h`, `network_data.c`, `network_config.h` and `network_config.h`) by those generated previously with STM32CubeMX. Replace them in Windows explorer. Below the file location:

```
PersonDetection +-- MobileNetv2_Model | +-- Inc | | network.h | | network_data.h | | network_config.h | +-- Src | | network.c | | network_data.c
```

 Change the name of your classes. Go to `fp_vision_app.c` and change `output_labels` to your class names. Be careful that there are as many labels as your model output size. Clean your project under STM32CubeIDE with `Project->Clean...`. Make sure that your board is connected to the ST-Link port and properly powered (set the jumper position to `ST-Link`



board is connected to the ST-LINK port and properly powered (set the jumper position to `ST-LINK`). Build your project and run it with the play button. A ST-Application pack splash window appears followed by the camera view and the model classification. If either of these window is not displayed properly, go through the wiki article to check if steps have been skipped. `<noinclude> [[Category:Artificial Intelligence|11]]`
`{{PublicationRequestId | 20512| 2021-07-12 | AnneJ ?}}` `</noinclude>`

Template used on this page:

- [Template:PublicationRequestId \(view source\)](#)

[Return to Main Page.](#)