



***NMDL - Software for the  
implementation of deep neural  
networks on the NeuroMatrix®  
platform***

## **NMDL User's manual**



# Contents

1. General information .....	6
1.1. The Purpose of Software .....	6
1.2. Quick start .....	6
1.3. NMDL performance .....	7
2. Software components .....	9
3. Installation .....	11
3.1. Windows installation .....	11
3.2. Linux installation .....	11
3.3. Additional components .....	11
3.4. Preparing the MC121.01 module .....	12
3.5. Preparing the NMStick device .....	13
3.6. Preparation of MC127.05 and NMCard modules .....	13
4. Compiling the model .....	14
4.1. Supported operations .....	14
5. Preparing images .....	17
6. Processing modes .....	19
7. Demo program .....	22
8. An example of using NMDL .....	28
9. Description of identifiers, functions and structures of NMDL .....	33
9.1. Identifiers and structures .....	33
9.1.1. NMDL_BOARD_TYPE .....	33
9.1.2. NMDL_ModelInfo .....	33
9.1.3. NMDL_PROCESS_FRAME_STATUS .....	33
9.1.4. NMDL_RESULT .....	34
9.1.5. NMDL_Tensor .....	35
9.2. Functions .....	35
9.2.1. NMDL_Blink .....	35
9.2.2. NMDL_Create .....	36
9.2.3. NMDL_Destroy .....	36
9.2.4. NMDL_GetBoardCount .....	36
9.2.5. NMDL_GetLibVersion .....	36
9.2.6. NMDL_GetModelInfo .....	37
9.2.7. NMDL_GetOutput .....	37
9.2.8. NMDL_GetStatus .....	37
9.2.9. NMDL_Initialize .....	38
9.2.10. NMDL_ProcessFrame .....	38
9.2.11. NMDL_Release .....	39
9.2.12. NMDL_RequiredOutputFloats .....	39
10. Description of identifiers and functions nmdl_compiler .....	40
10.1. Identifiers .....	40
10.1.1. NMDL_COMPILER_BOARD_TYPE .....	40
10.1.2. NMDL_COMPILER_RESULT .....	40
10.2. Functions .....	41
10.2.1. NMDL_COMPILER_CompilerDarkNet .....	41
10.2.2. NMDL_COMPILER_CompilerONNX .....	41

10.2.3. NMDL_COMPILER_FreeModel .....	42
10.2.4. NMDL_COMPILER_GetLastError .....	42
11. Description of identifiers and functions nmdl_image_converter .....	43
11.1. Identifiers and structures .....	43
11.1.1. NMDL_IMAGE_CONVERTER_BOARD_TYPE .....	43
11.1.2. NMDL_IMAGE_CONVERTER_COLOR_FORMAT .....	43
11.2. Functions .....	43
11.2.1. NMDL_IMAGE_CONVERTER_Convert .....	43
11.2.2. NMDL_IMAGE_CONVERTER_RequiredSize .....	44

## List of Figures

3.1. Jumpers on the MC121.01 module .....	13
6.1. NM6408 clusters. ....	19
6.2. Processing modes .....	20
7.1. The appearance of the program in the process .....	23
7.2. Module selection window appearance .....	24
7.3. Classification results window appearance .....	26
7.4. Displaying detection results .....	27

## List of Tables

1.1. FPS .....	7
1.2. Latency (ms) .....	8

# 1. General information

## 1.1. The Purpose of Software

*NMDL* software module allows you to run a pretrained deep convolutional neural network on computational modules *MC121.01*, *MC127.05*, *NMStick*, *NMCard* and on the *MC127.05* module simulator . The software module consists of 2 parts. One part runs on a personal computer (host) running 64-bit Microsoft® Windows 7/10 or Linux OS, the other part starts and runs on the processor of the computing module. Communication of *MC121.01* and *NMStick* devices with the host is carried out via the USB2.0 channel, for communication of *MC127.05* and *NMCard* modules with the host, the PCIe interface is used.

To work with *NMDL*, you must first install the software to support the computational modules used in the system. No support software installation is required to work with the simulator.

*NMDL* performs processing of custom source images in accordance with the specified neural network model. Before processing, you need to prepare the model and image data.

The model is preliminarily prepared by a special compiler from *NMDL*. Source models can be in *ONNX* or *DarkNet* format. Not all the operations defined in *ONNX* are supported by the *NMDL* compiler. For a list of supported operations and other restrictions, see "[Supported Operations](#)" .

Images must also be pre-processed with a special image converter. Only prepared models and images can be loaded and processed on computational modules.

The library provides a C / C ++ programming interface.

Further in the text *NMDL* (in upper case) - designation of the software package, *nmdl* (in lower case) - files of the program module.

## 1.2. Quick start

This section provides step-by-step instructions for getting started quickly with the [demo](#) . More complete information about *NMDL* can be found in the corresponding sections of the manual. Here is a demonstration of image processing using a neural network to classify *SqueezeNet* on a board with NM6408 processor simulator.

- Install the *NMDL* distribution as described in [Installation section](#).
- Go to the *bin* directory in the *NMDL* installation directory (by default in Windows "c: \ Program Files \ Module \ NMDL", by default in Linux "/ opt / nmdl ") and run the *nmdl\_gui* demo program.
- Select the device using the *File - Open Board ...* menu. Make sure the simulator is selected in the dialog box.
- Select the model description using the *File - Open Description...* menu. In the file selection dialog, select the file *PATH\_TO\_NMDL/nmdl\_ref\_data/ squeezeNet\_imagenet/ description08.xml*.

- Select the image to be processed using the *File - Open Picture...* menu. In the file selection dialog, select the *PATH\_TO\_NMDL/nmdl\_ref\_data/squeezeenet\_imagenet/frame.bmp* file.
- Start processing using the *File - Run* menu. As a result of processing, a classification window will pop up with the calculated probabilities in decreasing order. The correct class for the selected image is "lakeside".

## 1.3. NMDL performance

The tables show the performance values (frames per second FPS) and the delay values from the start of frame processing until the result is obtained (Latency). The size of the processed image is indicated in brackets next to the network name.

**Table 1.1. FPS**

	<b>MC121.01</b>	<b>NMStick</b>	<b>MC127.05 NMCard</b>	<b>MC127.05 NMCard batch-mode*</b>
alexnet (227x227)	3,45	3,2	12,6	13
inception v3 (299x299)	0,63	0,6	12,8	20,3
inception v3 (512x512)	0,24	0,23	3,93	5,44
resnet 18 (224x224)	2,28	2,2	25	47
resnet 50 (224x224)	0,8	0,75	12,2	20,6
squeezeenet (224x224)	8,3	8	74,4	100
u-net (512x512)**	-	-	2	2
yolo v2 tiny (416x416)	1,16	1,1	21	30,4
yolo v3 (416x416)	0,1	0,09	3,7	4,5
yolo v3 tiny (416x416)	1,44	1,38	27,3	35,3

**Table 1.2. Latency (ms)**

	<b>MC121.01</b>	<b>NMStick</b>	<b>MC127.05 NMCard</b>	<b>MC127.05 NMCard batch-mode*</b>
alexnet (227x227)	290	302	79	308
inception v3 (299x299)	1587	1653	78	197
inception v3 (512x512)	4166	4340	254	735
resnet 18 (224x224)	439	457	40	85
resnet 50 (224x224)	1250	1300	82	194
squeezezenet (224x224)	120	125	13	40
u-net (512x512)**	-	-	500	2000
yolo v2 tiny (416x416)	862	898	47	132
yolo v3 tiny (416x416)	10000	10416	270	889
yolo v3 tiny (416x416)	694	725	36	113

\* For a description of *batch-mode*, see ["Processing Modes"](#)

\*\* The model u-net has replaced the *transposed\_convolution* layers with *upsampling*.

## 2. Software components

Neural network implementation software consists of program modules (API), utilities and manuals.

API files for developing programs using *NMDL*:

- *nmdl.dll/nmdl.so* - software module for inference a neural network. See ["Description of nmdl identifiers, functions and structures"](#).
- *nmdl.lib* - library for early binding of programs with *NMDL* in MSVC ++ environment.
- *nmdl.h* - header file with description of API structures and functions.
- *nmdl\_compiler.dll/nmdl\_compiler.so* - program module - model compiler *ONNX / DarkNet* to internal representation. See ["Description of nmdl\\_compiler identifiers and functions"](#)
- *nmdl\_compiler.lib* - library for early binding of the model compiler module in the MSVC ++ environment.
- *nmdl\_image\_converter.dll/nmdl\_image\_converter.so* - программный модуль для подготовки обрабатываемых изображений. См. ["Описание идентификаторов и функций nmdl\\_image\\_converter"](#)
- *nmdl\_compiler.h* - header file describing the structures and functions of the model compiler.
- *nmdl\_image\_converter.dll / nmdl\_image\_converter.so* - a program module for preparing processed images. See ["Description of nmdl\\_image\\_converter identifiers and functions"](#)
- *nmdl\_image\_converter.lib* - module for early binding of the MSVC++ image preparation module.
- *nmdl\_image\_converter.h* - header file describing structures and functions for preparing images.

Header files and early binding libraries are located in the *include* and *lib* directories of the *NMDL* directory.

Utilities:

- *nmdl\_compiler\_console* - command line utility for compiling models from *ONNX* and *DarkNet* formats into internal format for loading on computational modules. The *ONNX* model file usually has the extension *.pb* . The model in *DarkNet* format is saved in two files - with the extension *.cfg* and the extension *.weights* . The prepared model for the *MC121.01* and *NMStick* devices has the extension *.nm7* . The model for *MC127.05* and *NMCard* has the extension *.nm8* . See ["Compiling a Model"](#).
- *nmdl\_nmdl\_image\_converter\_console* - command line utility for preparing processed images. See [Preparing Images](#).

- *nmdl\_gui* - a windowing utility to demonstrate the functionality of the NMDL. See "["Demo"](#)

## 3. Installation

### 3.1. Windows installation

*NMDL* distributed as an installation program. Only 64-bit systems are supported.

To install the distribution, run the installer executable file with administrator rights. Follow the instructions of the installation wizard.

To work with program modules, they must be included in the "visibility" of the operating system. One solution is to create an environment variable, for example, named *NMDL*, which records the path to the directory with header and binaries, and adding the created variable to the PATH environment variable.

To use the *nmdl* module in C / C ++ programs, include the #include "nmdl.h" directive in the source file and link the program with the *NMDL* library. If you are using the MSVC++ development environment, include the *nmdl.lib* file for early binding. You can create and use an environment variable *NMDL* to set the path to the files *nmdl.h* and *nmdl.lib*. In the same way, you can connect the *nmdl\_compiler* and *nmdl\_image\_converter* modules.

### 3.2. Linux installation

It is distributed as a .deb package. Only 64-bit systems of the Debian family are supported.

Use the dpkg package manager to install.

For example:

```
dpkg -i NMDL.deb
```

### 3.3. Additional components

Additional components are available with the software module to test and demonstrate *NMDL*.

Components are distributed in archives:

- *nmdl\_ref\_data\_alexnet\_imagenet.zip* - archive for demonstration of the *ALEXNET*.
- *nmdl\_ref\_data\_inception\_v3\_imagenet.zip* -archive for demonstration of the *INCEPTION V3*.
- *nmdl\_ref\_data\_resnet\_18\_imagenet.zip* -archive for demonstration of the *RESNET18*.
- *nmdl\_ref\_data\_resnet\_50\_imagenet.zip* -archive for demonstration of the *RESNET50*.
- *nmdl\_ref\_data\_squeezezenet\_imagenet.zip* - archive for demonstration of the *SQUEEZENET*.
- *nmdl\_ref\_data\_unet\_no\_transp\_conv\_covid.zip* - archive for demonstration of the *U-NET*. The model has replaced the *transposed\_convolution* layers with *upsampling*.

- *nmdl\_ref\_yolo2\_tiny\_pascal\_voc.zip* - archive for demonstration of the *YOLONET V2 TINY*.
- *nmdl\_ref\_yolo3\_coco.zip* - archive for demonstration of the *YOLONET V3*.
- *nmdl\_ref\_yolo3\_tiny\_coco.zip* - archive for demonstration of the *YOLONET V3 TINY*.

The directories contain files:

- *frame.bmp* - test image.
- *model.cfg* - model in *DARKNET* format.
- *model.pb* - model in *ONNX* format.
- *model.weights* - model weights in *DARKNET* format.
- *description07.xml* - model description file for *MC121.01* and *NMStick* devices to run in the *nmdl\_gui* program.
- *description08.xml* - model description file for modules *MC127.05*, *NMCard* and a simulator to run in the *nmdl\_gui* program.

To prepare demo data that will be processed on the computational module, you need to unpack it into the *nmdl\_ref\_data* directory, compile models and prepare images as described in "[Compile Model](#)" and "[Prepare Images](#)" of this manual. For ease of compilation, the archive includes scripts *prepare.cmd* and *prepare.sh*. As a result of the script's work, files will appear in each directory:

- *frame07* - image prepared for processing on *MC121.01* and *NMStick*.
- *frame08* - image prepared for processing on *MC127.05* and *NMCard*.
- *model.nm7* - compiled model for uploading to *MC121.01* and *NMStick*.
- *model.nm8* - compiled model for uploading to *MC127.05* and *NMCard*.

## 3.4. Preparing the MC121.01 module

When using NMDL with *MC121.01* devices, set the following jumpers before connecting the module to the PC USB connector (see figure [3.1](#)):

- Pin 1-2 of connector X9.
- Pins 3-4 of connector X9.
- Pins 5-6 of connector X9.
- When using USB power, pins 1-2 of connector X11 (when using a power supply, open the contacts).

- Pins 1-2 of connector X18.

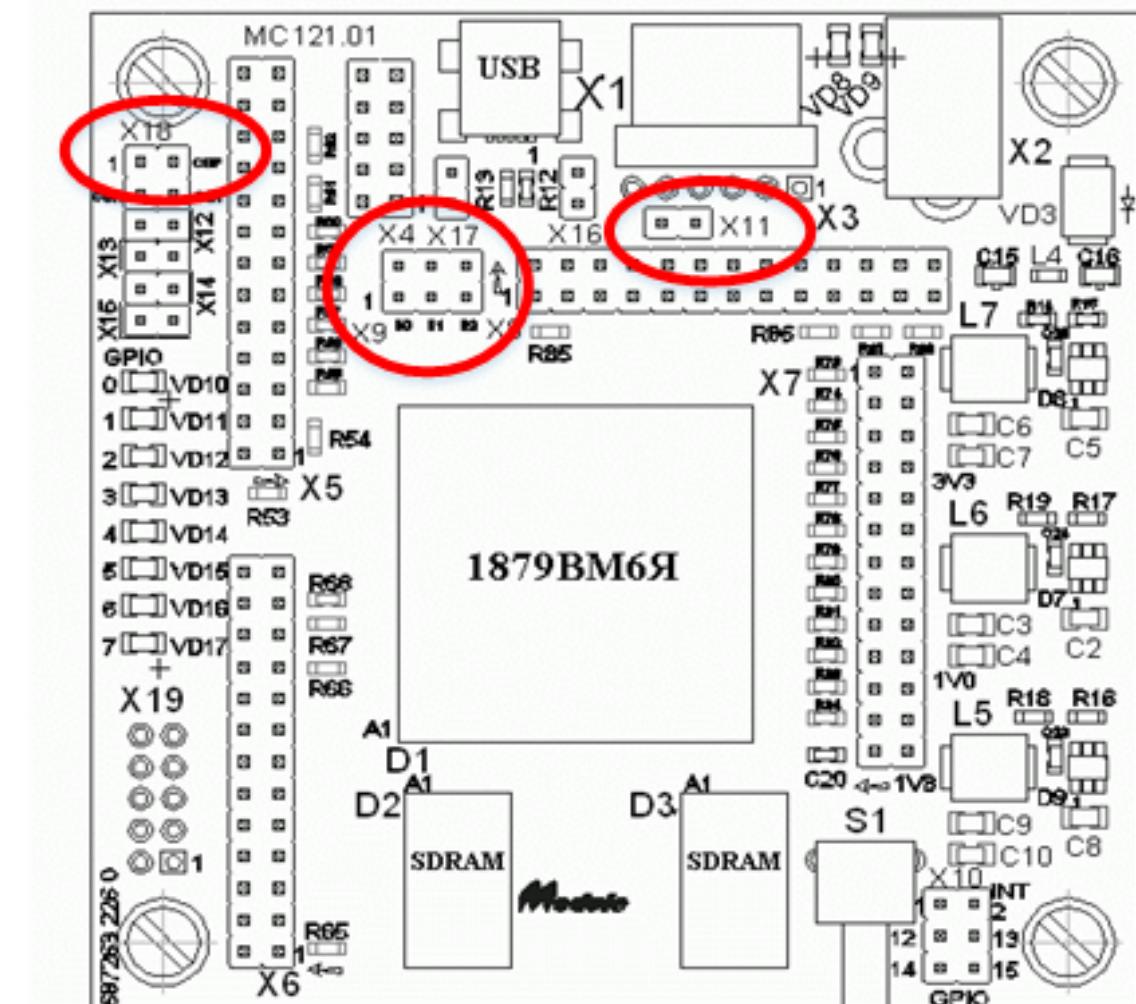


Figure 3.1. Jumpers on the MC121.01 module

### 3.5. Preparing the NMStick device

When using NMDL with *NMStick*, you need to install the device support software (supplied with the product) and connect the device to the USB port of your PC.

### 3.6. Preparation of MC127.05 and NMCard modules

To work with NMDL, you need to install the card in a free PCIe slot and install the module support software included in the product delivery set.

## 4. Compiling the model

The neural network must be pre-trained using a neural network package (e.g. Microsoft® Cognitive Toolkit (CNTK)), and converted to *ONNX* or *DarkNet*.

The original ONNX model is usually stored in a file with the Google Protocol Buffer structure and has a \* .pb or \* .onnx extension.

Models for YOLO networks can be stored in DarkNet format. In this case, the model is described by two files - a file with the \* .cfg extension contains a description of the processing graph, he \* .weights file contains the weights for the convolutions.

The ONNX and DarkNet model files are the result of training the neural network and at the same time the initial data for the compiler, which, as a result of processing, creates files with *nm7* extension for *MC121.01* and *NMStick* devices, or *nm8* for *MC127.05 modules*, *NMCard* and simulator.

The *nm7* and *nm8* models are binary images of compiled user models. Their structure is not documented and depends on the target device and compiler version.

The compiler is designed as a dynamically loadable module and can be embedded in a user program. You can also use the console utility `nmdl_compiler_console` to perform model transformations from the command line.

```
nmdl_compiler_console BOARD_TYPE NN_TYPE
    SRC_FILENAME DST_FILENAME [WEIGHTS_FILENAME]
```

Command line arguments:

- *BOARD\_TYPE* - module type. Valid values: "*MC12101*" - model conversion to nm7 format for *MC121.01* and *NMStick*, "*MC12705*" - converting the model to nm8 format for the *MC127.05*, *NMCard* modules and the simulator.
- *NN\_TYPE* - the file format of the input model. Valid values are "*ONNX*" or "*DARKNET*".
- *SRC\_FILENAME* - the filename of the original model.
- *DST\_FILENAME* - the filename of the output model.
- *WEIGHTS\_FILENAME* - the name of file with model weights. Needed only for models in DarkNet format.

Example:

```
>nmdl_compiler_console MC12705 ONNX
    squeezenet.pb squeezenet.nm8
>nmdl_compiler_console MC12705 DARKNET
    yolo2t.pb yolo2t.nm8 yolo2t.weights
```

### 4.1. Supported operations

The compiler supports the following set of operations:

1. Abs
2. Add
3. AveragePool
  - AveragePool2x2, no pad, stride=1
  - AveragePool2x2, no pad, stride=2
  - AveragePool3x3, no pad, stride=2
  - AveragePool3x3, pad, stride=2
4. BatchNormalization (provided that the operation is performed immediately after the convolution).
5. Concat (channels axis or width axis, the number of input tensors must be no more than 7)
6. Convolution
  - Conv1x1, stride=1
  - Conv1x1, stride=2
  - Conv3x3, no pad, all strides
  - Conv3x3, pad, stride=1
  - Conv3x3, pad, stride=2
  - Conv5x5, no pad, all strides
  - Conv5x5, pad, stride=1
  - Conv7x7, no pad, all strides
  - Conv7x7, pad, stride=2
  - Conv11x11, no pad, all strides
  - Conv7x1, pad\_w, stride=1
  - Conv1x7, pad\_h, stride=1
  - Conv3x1, pad\_w, stride=1
  - Conv1x3, pad\_h, stride=1
7. ConvTranspose (kernel 2x2, stride 2x2)
8. Div

9. GEMM
10. GlobalAveragePool
11. Leaky Relu
12. Mat Mul
13. MaxPool
  - MaxPool2x2, no pad, stride=1
  - MaxPool2x2, no pad, stride=2
  - MaxPool3x3, no pad, stride=2
  - MaxPool3x3, pad, stride=2
14. Mul
15. Pad
16. PRelu
17. Relu
18. Reshape
19. Sigmoid
20. Slice
21. Sub
22. Transpose

NMDL can work with models with one processed image, that is, one input tensor is processed.

The number of terminal nodes (output tensors) - no more than 5.

## 5. Preparing images

The processed images must first be converted to a float format (tensor). This can be done using the *nmdl\_image\_converter\_console* utility

```
nmdl_image_converter_console SRC_FILE DST_FILE W H F DR DG DB AR AG AB BT
```

Command line arguments:

- *SRC\_FILE* - input file name (*bmp*, *gif*, *jpg*, *emf*, *png*, *tiff*),
- *DST\_FILE* - output file name,
- *W* - image tensor width,
- *H* - image tensor height,
- *F* - the order of the color channels in the image tensor (valid values: *rgb* , *rbg* , *grb* , *gbr* , *brg* , *bgr* , intensity - one grayscale channel),
- *DR* - divisor for red pixel channel in the expression  $dst = src / D + A$  (float value),
- *DG* - divisor for green pixel channel in the expression  $dst = src / D + A$  (float value),
- *DB* - divisor for blue pixel channel in the expression  $dst = src / D + A$  (float value),
- *AR* - the term for red pixel channel in the expression  $dst = src / D + A$  (float value).
- *AG* - the term for green pixel channel in the expression  $dst = src / D + A$  (float value).
- *AB* - the term for blue pixel channel in the expression  $dst = src / D + A$  (float value).
- *BT* - the type of the module on which the processing is supposed (valid values: *mc12101*, *mc12705*).

For grayscale images, when F is intensity, only the divisor DR and the AR term are used. Other divisors (DG and DB) and terms (AG and AB) are ignored and can be set to any value.

The program scales the input image relative to the center according to the given arguments.

The prepared images are either planar (for *MC121.01* and *NMStick* ) or pixelated (for *MC127.05* and *NMCard* ) layout format for elements of type *float32* .

In the planar format, the planes of each channel are written to the file in turn. In this case, the buffer can be thought of as a C/C ++ style array:

```
float image[CHANNELS][HEIGHT][WIDTH];
```

That is, the fastest changing index is the index on the width of the image. The number of channels is three (RGB channels), or one for single-channel images (grayscale). In the prepared

---

buffer, the size is aligned to the width of the image. For images with even width, the buffer size will be:  $size = width * height * channels$  (float32), for images with odd width:  $size = (width + 1) * height * channels$  (float32).

In the pixel format, the channels of the image pixels are written to the file sequentially. A buffer can be thought of as a C / C ++ style array:

```
float image[HEIGHT][WIDTH][CHANNELS];
```

That is, the fastest changing index is the channel index. The number of channels is three (RGB channels), or one for single-channel images (grayscale). In the prepared buffer, the size is aligned along the image channels to the nearest even value. The size of the buffer with the prepared image for the *MC127.05* and *NMCard* modules will be:  $size = width * height * (channels + 1)$  (float32).

## 6. Processing modes

The architectural features used in the computational modules allow implementing neural network algorithms in different ways. The variety is determined by the hardware to do parallel computing. The implementation of *NMDL* used deep optimization of algorithms and chose the fastest parallelism methods, encapsulating development details. Nevertheless, in the implementation of *NMDL* for the MC127.05 module, there remains a variability of solutions, which it is advisable to put into the interface and provide the user with the opportunity to choose one or another processing method depending on the problem being solved.

Used in the modules *MC127.05* and *NMCard* NM6408 chip has a cluster organization.

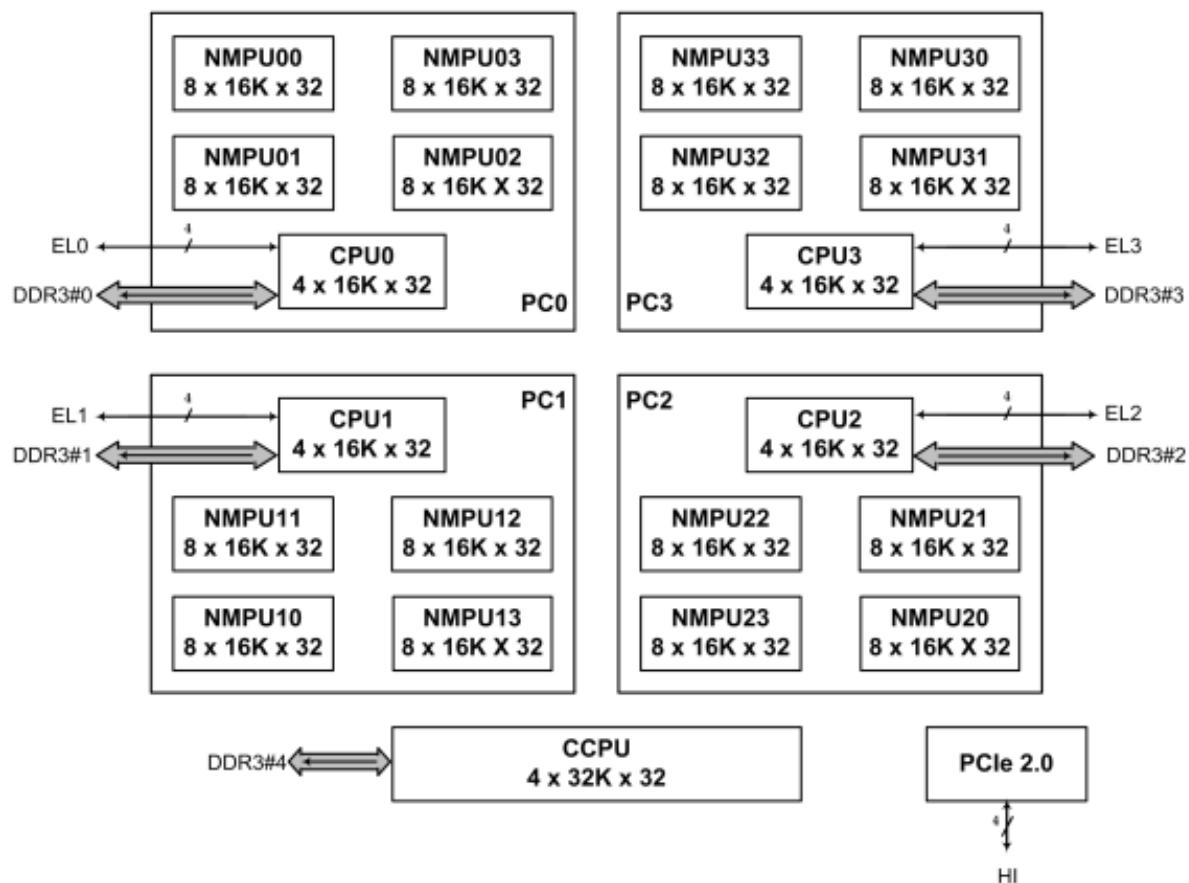


Figure 6.1. NM6408 clusters.

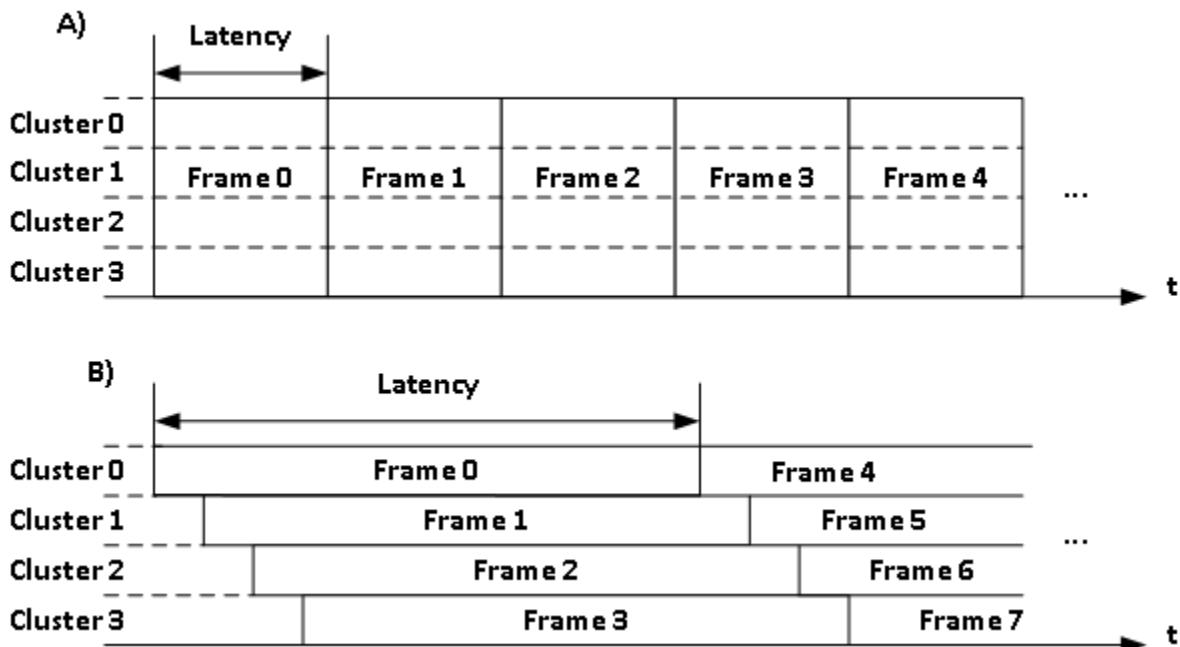
On the picture [6.1](#):

- *PC* - processor clusters;

- *NMPU* - NMC4 processor nodes;
- *CPU* - cluster processor cores ARM Cortex-A5;
- *CCPU* - central processor core ARM Cortex-A5.

Within the clusters, the amount of computation and data is distributed over four processor nodes NMC4. Here, tight hardware connections are used for interprocessor communication within clusters and deep optimization is performed on primitive functions.

Processing on clusters can be organized in two ways: with data sharing between clusters and work in batch mode (*batch-mode*).



**Figure 6.2. Processing modes**

When processing with data sharing (Figure 6.2 A)), you can achieve the minimum latency - the time from the start of processing to the moment the result is obtained. However, in this case, parallelism overhead is inevitable. For example, when performing convolutions on split tensors, it is necessary to compensate for border processing, perform repackaging, data transit between clusters, etc. Performance - frames per second -  $FPS = 1 / Latency$ .

Batch processing (Figure 6.2 B)) processes one frame per cluster. The clusters perform the same and independent processing. In this case, there is no overhead for organizing data distribution and maximum performance is achieved, while latency increases. Performance - frames per second -  $FPS = 4 / Latency$ .

Changing of processing modes is possible only when working with the modules *MC127.05* , *NMCard* or with a simulator.

## 7. Demo program

To demonstrate the functionality of the NMDL library, the `nmdl_gui` utility was developed. The tasks of the utility include:

- Library initialization *NMDL*;
- Detection and identification of available accelerators (*simulator*, *MC121.01*, *MC127.05*, *NMStick*, *NMCard*);
- Loading a neural network model file (*.nm7* or *.nm8*);
- Loading the neural network description file (*xml*);
- Pre-processing and loading of images on the computer module;
- Starting processing on the computing module;
- Processing and output of results.

When the program is started, the main program window appears with menu bars, status and client area.

The menu contains controls for the program. The client area contains the image and processing results. The status bar displays the following information:

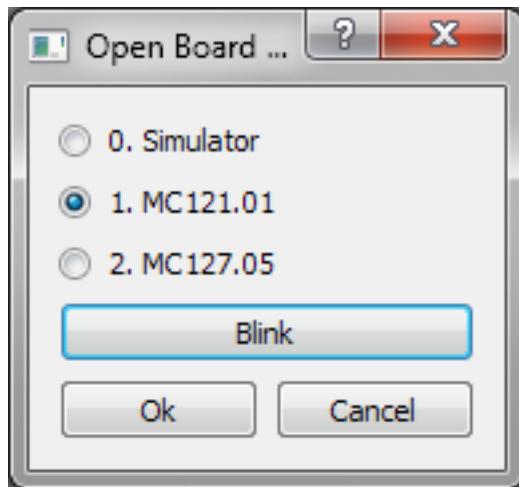
- Selected accelerator (*simulator*, *MC121.01*, *MC127.05*, *NMStick*, *NMCard*);
- Selected neural network model;
- Selected description of the neural network;
- The selected image;
- Processing speed (frames/s).

The appearance of the program during operation is shown in the figure [7.1](#).



**Figure 7.1. The appearance of the program in the process**

The choice of the module is carried out in the *Open Board* dialog box, the appearance of which is shown in the figure 7.2. The window contains a list of accelerators available in the system, selection buttons, and the ability to identify the device by means of LED indication.



**Figure 7.2. Module selection window appearance**

The selection of the neural network description is carried out in the *Open Description* dialog box. The neural network is described in *XML* format and contains information on the required image format, as well as parameters for interpreting the output parameters.

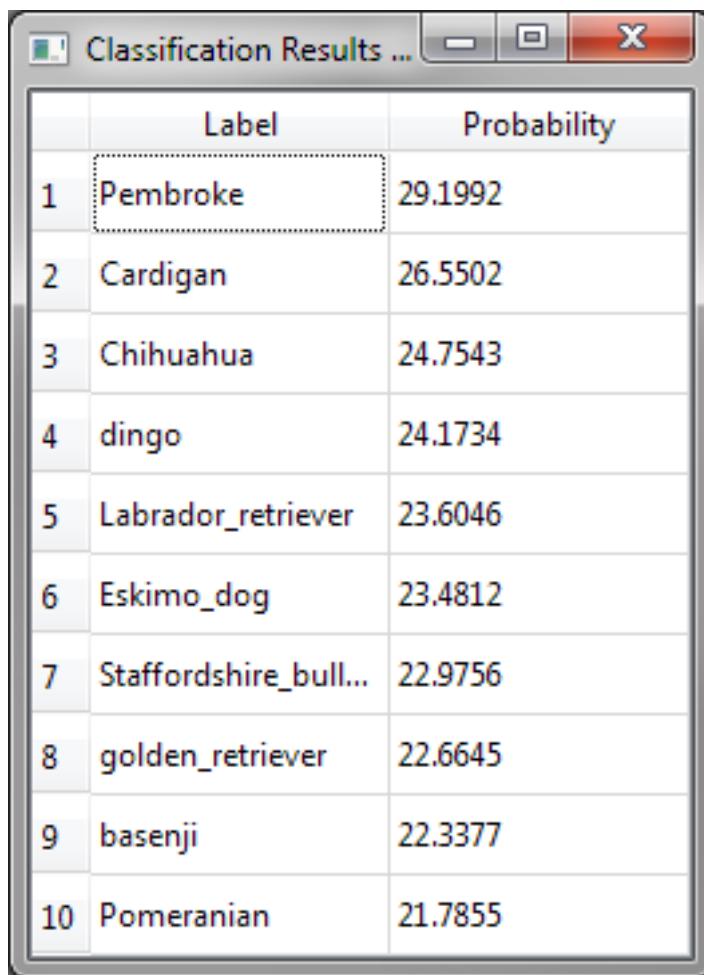
Parameter	Description
model	<p>A string with the name of the neural network model file in one of the following formats:</p> <ul style="list-style-type: none"> <li>• *.nm7 - description of the model for loading on the <i>MC121.01</i> and <i>NMStick</i> modules.</li> <li>• *.nm8 - description of the model for loading on the <i>MC127.05</i>, <i>NMCard</i> modules or on the simulator.</li> <li>• *.pb or *.onnx - description of the model in ONNX Protobuf format. The model can be loaded onto any module. Before loading onto a module, the model is pre-converted to *.nm7 or *.nm8 depending on the type of the module.</li> <li>• *.cfg - description of the model in DARKNET format. The file with weights *.weights must be located in the same directory and have the same name as the model description file *.cfg. The model can be loaded onto any module. Before loading onto a module, the model is pre-converted to *.nm7 or *.nm8 depending on the type of the module.</li> </ul> <p>You can specify either an absolute path to the file or a path relative to the location of the XML deployment descriptor.</p>

Parameter	Description
format	<p>The pixel format to which the pixels in the input image will be converted before processing. Can take values:</p> <ul style="list-style-type: none"> <li>• rgb</li> <li>• rbg</li> <li>• grb</li> <li>• gbr</li> <li>• brg</li> <li>• bgr</li> </ul> <p>This value corresponds to the order of the input tensor channels.</p>
divider	The scaling factor for each color component of a pixel in the input image. Used when converting an image to an input tensor (see <a href="#">"Preparing Images"</a> ).
adder	Displacement factor for each color component of a pixel in the input image. Used when converting an image to an input tensor (see <a href="#">"Preparing Images"</a> ).
neural_network	<p>Selected neural network Can take values:</p> <ul style="list-style-type: none"> <li>• <i>classifier</i> - neural network - a classifier, for example <i>ALEXNET</i>, <i>RESNET</i>, <i>SQUEEZENET</i>.</li> <li>• <i>unet</i> - neural network of the U-Net family.</li> <li>• <i>yolo2</i> - neural network of the YOLO V2 family.</li> <li>• <i>yolo3</i> - neural network of the YOLO V3 family.</li> </ul> <p>The result of the work of the classifier is a vector of probabilities that the image belongs to certain classes. In the demo, the classes and probabilities of detection are displayed in a pop-up window.</p> <p>As a result of processing, YOLO networks provide information about several detected objects and their location in the original image. In the demo program, the detected objects are indicated directly on the original image in the enclosing rectangles.</p>
yolo_anchors	Parameters for initial initialization of detected rectangles (only for <i>YOLO2</i> and <i>YOLO3</i> networks).
yolo_confidence_threshold	Threshold for displaying detection results (only for <i>YOLO2</i> and <i>YOLO3</i> networks).

Parameter	Description
yolo_iou_threshold	Overlapping detection rectangle threshold (only for <i>YOLO2</i> and <i>YOLO3</i> networks)
labels	List of strings with class names that will be displayed in the processing result window.

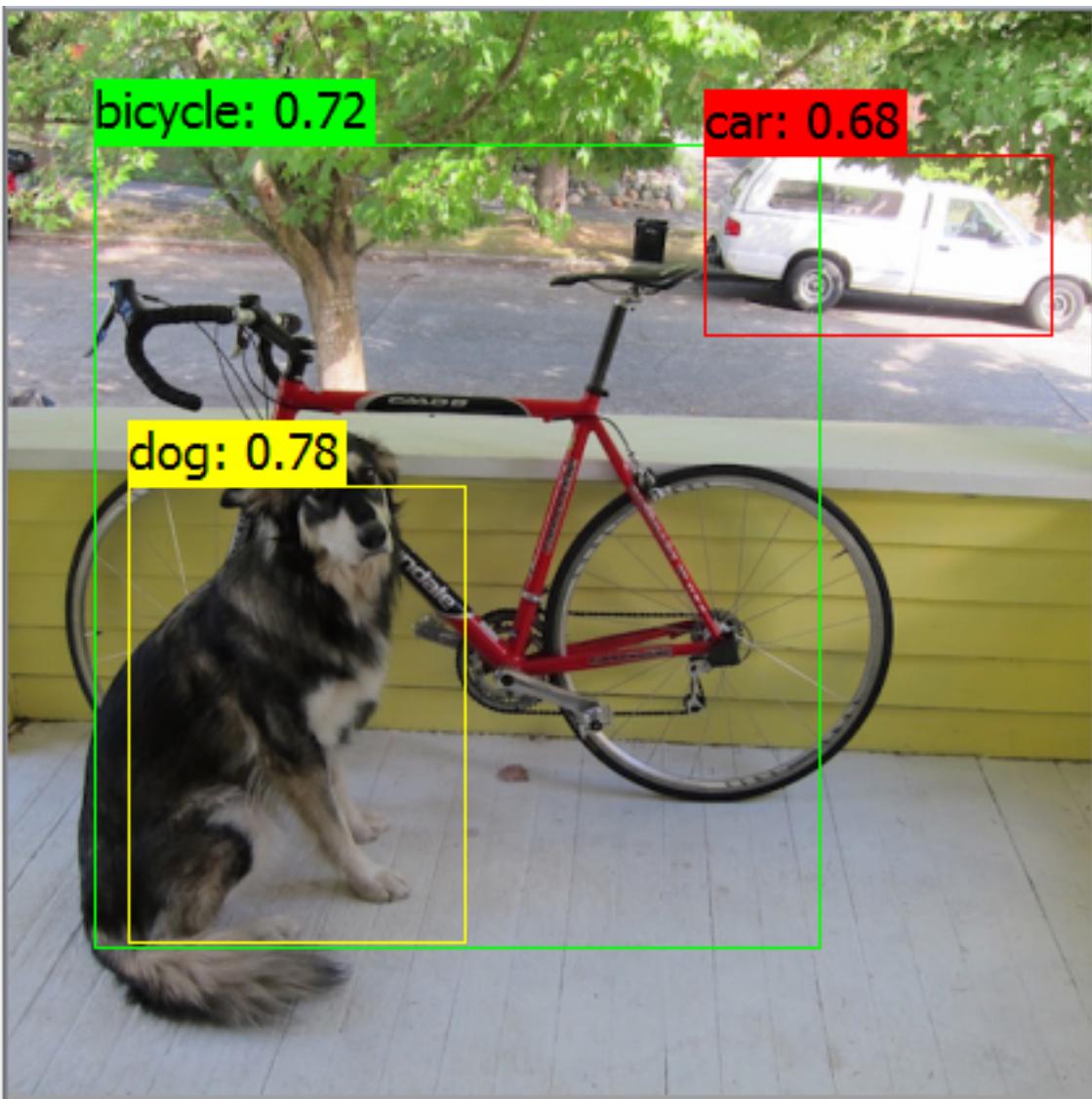
The *Open Picture* dialog box allows you to open one or more images for processing.

When all the data is available, the *Run* button becomes available and starts processing. The launch can be initiated using the "*Ctrl + R*" combination. Pre-processing is carried out for each selected image in accordance with the given description of the neural network for further processing on the accelerator. As a result of processing, the output structure with N tensors is filled. Depending on the neural network, the output structure is interpreted in different ways. For classification networks (such as *SqueezeNet*), an additional window with classes and probability appears. The appearance of this window is shown in the figure [7.3](#).



**Figure 7.3. Classification results window appearance**

For networks that perform object detection (such as *YOLO*), the result is shown directly in the original image. After the accelerator is finished, the rectangles are superimposed on the detected objects, the class name and the detection probability are given. An example of overlaying the result on an image is shown in the figure [7.4](#).



**Figure 7.4. Displaying detection results**

The program allows you to process a sequence of images. To do this, select several files. After that, when starting processing (*Run*), one image will be processed. When restarting - the second, etc. You can also start processing in automatic mode (*Run Auto*). Press the *spacebar* key to stop automatic processing.

The status bar is used to display additional information. The line displays the name of the selected device, the neural network description file name, the image filename, as well as the fps value measured by the program, taking into account the frame transfer and the result obtained. The processing speed without taking into account the transfer time is displayed in brackets.

## 8. An example of using NMDL

The example demonstrates the use of the *NMDL* library, software modules for model compilation and image preparation. The example consists of a source file *example.cpp* and a build script *CMakeLists.txt* in the "*examples*" directory of the installation directory.

It is assumed that the source data for processing is in the "*nmdl\_ref\_data/squeezeenet*" folder in the installation directory, this means that the example demonstrates the processing of the neural network *squeezeenet*. The initial data are:

- Neural network model in ONNX format - file "*nmdl\_ref\_data/squeezeenet/model.pb*"
- Processed image in BMP format - file "*nmdl\_ref\_data/squeezeenet/frame.bmp*"

The example shows the calls to the functions of the libraries in the order necessary for correct operation.

```

001 #include <fstream>
002 #include <iostream>
003 #include <string>
004 #include <unordered_map>
005 #include <vector>
006 #include "nmdl.h"
007 #include "nmdl_compiler.h"
008 #include "nmdl_image_converter.h"
009
010 namespace {
011
012 auto Call(NMDL_COMPILER_RESULT result, const std::string &function_name) {
013     static std::unordered_map<NMDL_COMPILER_RESULT, std::string> map = {
014         {NMDL_COMPILER_RESULT_OK, "OK"},
015         {NMDL_COMPILER_RESULT_MEMORY_ALLOCATION_ERROR, "MEMORY_ALLOCATION_ERROR"},  

016         {NMDL_COMPILER_RESULT_MODEL_LOADING_ERROR, "MODEL_LOADING_ERROR"},  

017         {NMDL_COMPILER_RESULT_INVALID_PARAMETER, "INVALID_PARAMETER"},  

018         {NMDL_COMPILER_RESULT_INVALID_MODEL, "INVALID_MODEL"},  

019         {NMDL_COMPILER_RESULT_UNSUPPORTED_OPERATION, "UNSUPPORTED_OPERATION"}  

020     };
021     if(result != NMDL_RESULT_OK) {
022         throw std::runtime_error(function_name + ": " + map[result] + ": " +
023             NMDL_COMPILER_GetLastError());
024     }
025     return NMDL_RESULT_OK;
026 }
027
028 auto Call(NMDL_RESULT result, const std::string &function_name) {
029     static std::unordered_map<NMDL_RESULT, std::string> map = {
030         {NMDL_RESULT_OK, "OK"},  

031         {NMDL_RESULT_INVALID_FUNC_PARAMETER, "INVALID_FUNC_PARAMETER"},  

032         {NMDL_RESULT_NO_LOAD_LIBRARY, "NO_LOAD_LIBRARY"},  

033         {NMDL_RESULT_NO_BOARD, "NO_BOARD"},  

034         {NMDL_RESULT_BOARD_RESET_ERROR, "BOARD_RESET_ERROR"},  

035         {NMDL_RESULT_INIT_CODE_LOADING_ERROR, "INIT_CODE_LOADING_ERROR"},  

036         {NMDL_RESULT_CORE_HANDLE_RETRIEVAL_ERROR, "CORE_HANDLE_RETRIEVAL_ERROR"},  

037         {NMDL_RESULT_FILE_LOADING_ERROR, "FILE_LOADING_ERROR"},  

038         {NMDL_RESULT_MEMORY_WRITE_ERROR, "MEMORY_WRITE_ERROR"},  

039         {NMDL_RESULT_MEMORY_READ_ERROR, "MEMORY_READ_ERROR"},  

040         {NMDL_RESULT_MEMORY_ALLOCATION_ERROR, "MEMORY_ALLOCATION_ERROR"},  

041         {NMDL_RESULT_MODEL_LOADING_ERROR, "MODEL_LOADING_ERROR"},  

042         {NMDL_RESULT_INVALID_MODEL, "INVALID_MODEL"},  

043         {NMDL_RESULT_BOARD_SYNC_ERROR, "BOARD_SYNC_ERROR"},  

044         {NMDL_RESULT_BOARD_MEMORY_ALLOCATION_ERROR, "BOARD_MEMORY_ALLOCATION_ERROR"},  

045         {NMDL_RESULT_NN_CREATION_ERROR, "NN_CREATION_ERROR"},  

046         {NMDL_RESULT_NN_LOADING_ERROR, "NN_LOADING_ERROR"},  

047         {NMDL_RESULT_NN_INFO_RETRIEVAL_ERROR, "NN_INFO_RETRIEVAL_ERROR"},  


```

```

048 {NMDL_RESULT_MODEL_IS_TOO_BIG,           "MODEL_IS_TOO_BIG"},  

049 {NMDL_RESULT_NOT_INITIALIZED,            "NOT_INITIALIZED"},  

050 {NMDL_RESULT_BUSY,                      "BUSY"},  

051 {NMDL_RESULT_UNKNOWN_ERROR,             "UNKNOWN_ERROR"}  

052 };  

053 if(result != NMDL_RESULT_OK) {  

054     throw std::runtime_error(function_name + ": " + map[result]);  

055 }  

056     return NMDL_RESULT_OK;  

057 }  

058  

059 template <typename T>  

060 auto ReadFile(const std::string &filename) {  

061     std::ifstream ifs(filename, std::ios::binary | std::ios::ate);  

062     if(!ifs.is_open()) {  

063         throw std::runtime_error("Unable to open input file: " + filename);  

064     }  

065     auto fsize = static_cast<std::size_t>(ifs.tellg());  

066     ifs.seekg(0);  

067     std::vector<T> data(fsize / sizeof(T));  

068     ifs.read(reinterpret_cast<char*>(data.data()), data.size() * sizeof(T));  

069     return data;  

070 }  

071  

072 void ShowNMDLVersion() {  

073     std::uint32_t major = 0;  

074     std::uint32_t minor = 0, patch = 0;  

075     Call(NMDL_GetLibVersion(&major, &minor, &patch), "GetLibVersion");  

076     std::cout << "Lib version: " << major << "." << minor << "." << patch << std::endl;  

077 }  

078  

079 void CheckBoard(std::uint32_t required_board_type) {  

080     std::uint32_t boards;  

081     std::uint32_t board_number = -1;  

082     Call(NMDL_GetBoardCount(required_board_type, &boards), "GetBoardCount");  

083     std::cout << "Detected boards: " << boards << std::endl;  

084     if(!boards) {  

085         throw std::runtime_error("Board not found");  

086     }  

087 }  

088  

089 auto CompileModel(const std::string &filename, std::uint32_t board_type) {  

090     auto onnx_model = ReadFile<char>(filename);  

091     float *nm_model = nullptr;  

092     std::uint32_t nm_model_floats = 0;  

093     Call(NMDL_COMPILER_CompileONNX(board_type, onnx_model.data(),  

094         onnx_model.size(), &nm_model, &nm_model_floats), "CompileONNX");  

095     std::vector<float> result(nm_model, nm_model + nm_model_floats);  

096     NMDL_COMPILER_FreeModel(board_type, nm_model);  

097     return result;  

098 }  

099  

100 auto GetModelInformation(NMDL_HANDLE nmdl) {  

101     NMDL_ModelInfo model_info;  

102     Call(NMDL_GetModelInfo(nmdl, &model_info), "GetModelInfo");  

103     std::cout << "Output tensor number: " << model_info.output_tensor_num << std::endl;  

104     for(std::size_t i = 0; i < model_info.output_tensor_num; ++i) {  

105         std::cout << "Output tensor " << i << ": " <<  

106         model_info.output_tensors[0].width << ", " <<  

107         model_info.output_tensors[0].height << ", " <<  

108         model_info.output_tensors[0].depth <<  

109         std::endl;  

110     }  

111     return model_info;  

112 }  

113  

114 auto PrepareFrame(const std::string &filename, std::uint32_t width,  

115     std::uint32_t height, std::uint32_t board_type,  

116     std::uint32_t color_format, float divider, float adder) {  

117     auto bmp_frame = ReadFile<char>(filename);  

118     std::vector<float> nm_frame(NMDL_IMAGE_CONVERTER_RequiredSize(

```

```

119     width, height, board_type));
120     if(NMDL_IMAGE_CONVERTER_Convert(bmp_frame.data(), nm_frame.data(), bmp_frame.size(),
121         width, height, color_format, divider, adder, board_type)) {
122         throw std::runtime_error("Image conversion error");
123     }
124     return nm_frame;
125 }
126
127 void WaitForOutput(NMDL_HANDLE nmdl, std::uint32_t batch_num,
128     std::vector<float> &output) {
129     std::uint32_t status = NMDL_PROCESS_FRAME_STATUS_BUSY;
130     while(status == NMDL_PROCESS_FRAME_STATUS_BUSY) {
131         NMDL_GetStatus(nmdl, batch_num, &status);
132     };
133     double fps;
134     Call(NMDL_GetOutput(nmdl, batch_num, output.data(), &fps), "GetOutput");
135     if(batch_num == 0) {
136         std::cout << "FPS: " << fps << std::endl;
137     }
138     std::cout << "Data: " << std::endl;
139     for(auto &cur: output) {
140         std::cout << "\t" << cur << std::endl;
141     }
142 }
143
144 }
145
146 int main() {
147     const std::uint32_t BOARD_TYPE = NMDL_BOARD_TYPE_SIMULATOR;
148 //const uint32_t BOARD_TYPE = NMDL_BOARD_TYPE_MC12705;
149 //const uint32_t BOARD_TYPE = NMDL_BOARD_TYPE_MC12101;
150     const std::uint32_t COMPILER_BOARD_TYPE =
151     BOARD_TYPE == NMDL_BOARD_TYPE_MC12101 ?
152     NMDL_COMPILER_BOARD_TYPE_MC12101 :
153     NMDL_COMPILER_BOARD_TYPE_MC12705;
154     const std::uint32_t IMAGE_CONVERTER_BOARD_TYPE =
155     BOARD_TYPE == NMDL_BOARD_TYPE_MC12101 ?
156     NMDL_IMAGE_CONVERTER_BOARD_TYPE_MC12101 :
157     NMDL_IMAGE_CONVERTER_BOARD_TYPE_MC12705;
158 #ifdef _WIN32
159     const std::string ONNX_MODEL_FILENAME = "..\\nmdl_ref_data\\squeezenet\\model.pb";
160     const std::string BMP_FRAME_FILENAME = "..\\nmdl_ref_data\\squeezenet\\frame.bmp";
161 #else
162     const std::string ONNX_MODEL_FILENAME = "../nmdl_ref_data/squeezenet/model.pb";
163     const std::string BMP_FRAME_FILENAME = "../nmdl_ref_data/squeezenet/frame.bmp";
164 #endif
165
166     const NMDL_IMAGE_CONVERTER_COLOR_FORMAT IMAGE_CONVERTER_COLOR_FORMAT =
167     NMDL_IMAGE_CONVERTER_COLOR_FORMAT_BGR;
168     const float NM_FRAME_DIVIDER = 1.0f;
169     const float NM_FRAME_ADDER = -114.0f;
170
171     const std::size_t BATCHES = 4;
172     const std::size_t FRAMES = 5;
173
174     NMDL_HANDLE nmdl = 0;
175
176     try {
177         std::cout << "Query library version..." << std::endl;
178         ShowNMDLVersion();
179
180         std::cout << "Board detection..." << std::endl;
181         CheckBoard(BOARD_TYPE);
182
183         std::cout << "Compile model..." << std::endl;
184         auto nm_model = CompileModel(ONNX_MODEL_FILENAME, COMPILER_BOARD_TYPE);
185
186         std::cout << "NMDL initialization..." << std::endl;
187         Call(NMDL_Create(&nmdl), "Create");
188         Call(NMDL_Initialize(nmdl, BOARD_TYPE, 0, nm_model.data(),
189             nm_model.size(), 0), "Initialize");

```

```

190
191     std::cout << "Get model information... " << std::endl;
192     auto model_info = GetModelInformation(nmdl);
193
194     std::cout << "Prepare frame... " << std::endl;
195     auto frame = PrepareFrame(BMP_FRAME_FILENAME, model_info.input_tensor.width,
196                               model_info.input_tensor.height, IMAGE_CONVERTER_BOARD_TYPE,
197                               IMAGE_CONVERTER_COLOR_FORMAT, NM_FRAME_DIVIDER, NM_FRAME_ADDER);
198
199     std::vector<float> output(NMDL_RequiredOutputFloats(&model_info));
200
201     std::cout << "Process single frames... " << std::endl;
202     for(std::size_t i = 0; i < FRAMES; ++i) {
203         Call(NMDL_ProcessFrame(nmdl, 0, frame.data()), "ProcessFrame");
204         WaitForOutput(nmdl, 0, output);
205     }
206     NMDL_Release(nmdl);
207
208     std::cout << "Batch mode process frames... " << std::endl;
209     std::uint32_t cnt_in = 0;
210     std::uint32_t cnt_out = 0;
211     Call(NMDL_Initialize(nmdl, BOARD_TYPE, 0, nm_model.data(),
212                           nm_model.size(), 1), "Initialize");
213     for(auto i = 0u; i < BATCHES; ++i) {
214         Call(NMDL_ProcessFrame(nmdl, (cnt_in++) % BATCHES,
215                               frame.data()), "ProcessFrame");
216     }
217     for(auto i = BATCHES; i < FRAMES; ++i) {
218         WaitForOutput(nmdl, (cnt_out++) % BATCHES, output);
219         Call(NMDL_ProcessFrame(nmdl, (cnt_in++) % BATCHES,
220                               frame.data()), "ProcessFrame");
221     }
222     for(auto i = 0u; i < BATCHES; ++i) {
223         WaitForOutput(nmdl, (cnt_out++) % BATCHES, output);
224     }
225 }
226 catch (std::exception& e) {
227     std::cerr << e.what() << std::endl;
228 }
229 NMDL_Release(nmdl);
230 NMDL_Destroy(nmdl);
231
232     return 0;
233 }
```

The following steps are performed here:

*1 .. 8*: Connecting header files. "*nmdl.h*" - description of the neural network processing library, "*nmdl\_compiler.h*" - description of the library for compiling models, "*nmdl\_image\_converter.h*" is a description of the library for preparing images.

*12 .. 26*: A wrapper function for calls to model compilation library functions. In case of an error, it generates an error message and raises an exception.

*28 .. 57*: A wrapper function for calls to functions of the neural network processing library. In case of an error, it generates an error message and raises an exception.

*59 .. 70*: Generic function for reading data from a file into a vector.

*72 .. 77*: Function for outputting the NMDL version [NMDL\\_GetLibVersion](#).

*79 .. 87*: Function for checking the presence of a module of a given type. In fact, a request is made for the number of detected modules of a given type - [NMDL\\_GetBoardCount](#).

---

89 .. 98 : Source model compilation function. The function is called [NMDL\\_COMPILER\\_CompilerONNX](#). This is where memory is allocated inside the compilation function, so after work the allocated memory is freed by calling [NMDL\\_COMPILER\\_FreeModel](#), and the compilation result is copied into the returned float vector. In the target program, you can pre-compile the model using the *nmdl\_compiler\_console* utility like this, as described in ["Compiling the Model"](#).

100..112: Function for receiving and displaying information about the parameters of the output tensors. The call is [NMDL\\_GetModelInfo](#).

114 .. 125: Frame preparation function. This is where you read an image from a file, decode it and preprocess it. ([NMDL\\_IMAGE\\_CONVERTER\\_Convert](#)). For a description of preprocessing, see [Preparing Images](#).

127 .. 142: Frame processing wait function. Takes the number of the cluster on which processing is performed (*batch\_num*) and vector buffer for storing the processing result. The function is blocked in the processing status request cycle ([NMDL\\_GetStatus](#)). After receiving the *NMDL\_PROCESS\_FRAME\_STATUS\_FREE* status, the result is copied by calling [NMDL\\_GetOutput](#).

158 .. 164: File names with the original model and image are specified.

166 .. 169: Parameters for preparing the image are set. For the original model, you need to prepare an image with pixels in the blue-green-red format. Each pixel is modified with the formula  $Y = X / 1.0 - 114$ . This transformation is required to process the *squeezeNet* model. For other models, it is necessary to use the corresponding parameters that are determined during model creation and are the initial data bound to a specific model. For more information about preparing images, see ["Preparing Images"](#).

171: Sets the number of clusters for batch processing.

172: The number of frames to be processed is set. In the example, one frame is processed multiple times.

Further, in the main function, the described auxiliary functions are sequentially called.

187 .. 189: Initializing NMDL. By calling the [NMDL\\_Create](#) function, the internal structures of the NMDL library are initialized. The [NMDL\\_Initialize](#) function loads the compiled model into the selected accelerator.

201 .. 205: An example of sequential processing of frames in the mode of dividing data into clusters (see the section ["Processing modes"](#)).

208 .. 224: An example of sequential processing of frames in batch processing mode (see the section ["Processing modes"](#)).

229 .. 230: Completion releases the allocated resources ([NMDL\\_Release](#) and [NMDL\\_Destroy](#)).

# 9. Description of identifiers, functions and structures of NMDL

## 9.1. Identifiers and structures

### 9.1.1. NMDL\_BOARD\_TYPE

Module types.

```
typedef enum tagNMDL_BOARD_TYPE {
    NMDL_BOARD_TYPE_SIMULATOR,
    NMDL_BOARD_TYPE_MC12101,
    NMDL_BOARD_TYPE_MC12705,
    NMDL_BOARD_TYPE_NMSTICK,
    NMDL_BOARD_TYPE_NMCARD
} NMDL_BOARD_TYPE;
```

- *NMDL\_BOARD\_TYPE\_SIMULATOR* - MC127.05 simulator.
- *NMDL\_BOARD\_TYPE\_MC12101* - MC121.01.
- *NMDL\_BOARD\_TYPE\_MC12705* - MC127.05.
- *NMDL\_BOARD\_TYPE\_NMSTICK* - NMStick.
- *NMDL\_BOARD\_TYPE\_NMCARD* - NMCard.

### 9.1.2. NMDL\_ModelInfo

Structure with information about the model.

```
typedef struct tagNMDL_ModelInfo {
    NMDL_Tensor input_tensor;
    unsigned int output_tensor_num;
    NMDL_Tensor output_tensors[NMDL_MAX_OUTPUT_TENSORS];
} NMDL_ModelInfo;
```

- *input\_tensor* - description of the input tensor (input image),
- *output\_tensor\_num* - number of output tensors,
- *output\_tensors* – output tensors.

*NMDL\_MAX\_OUTPUT\_TENSORS* - maximum number of output tensors.

### 9.1.3. NMDL\_PROCESS\_FRAME\_STATUS

Frame processing status identifier.

```
typedef enum tagNMDL_PROCESS_FRAME_STATUS {
```

```

NMDL_PROCESS_FRAME_STATUS_FREE,
NMDL_PROCESS_FRAME_STATUS_BUSY
} NMDL_PROCESS_FRAME_STATUS;

```

- *NMDL\_PROCESS\_FRAME\_STATUS\_FREE* - the frame is not processed,
- *NMDL\_PROCESS\_FRAME\_STATUS\_BUSY* - the frame is being processed.

## 9.1.4. NMDL\_RESULT

Returned result.

```

typedef enum tagNMDL_RESULT {
    NMDL_RESULT_OK,
    NMDL_RESULT_INVALID_FUNC_PARAMETER,
    NMDL_RESULT_NO_BOARD,
    NMDL_RESULT_BOARD_RESET_ERROR,
    NMDL_RESULT_INIT_CODE_LOADING_ERROR,
    NMDL_RESULT_CORE_HANDLE_RETRIEVAL_ERROR,
    NMDL_RESULT_FILE_LOADING_ERROR,
    NMDL_RESULT_MEMORY_WRITE_ERROR,
    NMDL_RESULT_MEMORY_READ_ERROR,
    NMDL_RESULT_MEMORY_ALLOCATION_ERROR,
    NMDL_RESULT_MODEL_LOADING_ERROR,
    NMDL_RESULT_INVALID_MODEL,
    NMDL_RESULT_BOARD_SYNC_ERROR,
    NMDL_RESULT_BOARD_MEMORY_ALLOCATION_ERROR,
    NMDL_RESULT_NN_CREATION_ERROR,
    NMDL_RESULT_NN_LOADING_ERROR,
    NMDL_RESULT_NN_INFO_RETRIEVAL_ERROR,
    NMDL_RESULT_MODEL_IS_TOO_BIG,
    NMDL_RESULT_NOT_INITIALIZED,
    NMDL_RESULT_BUSY,
    NMDL_RESULT_UNKNOWN_ERROR
} NMDL_RESULT;

```

- *NMDL\_RESULT\_OK* - no errors,
- *NMDL\_RESULT\_INVALID\_FUNC\_PARAMETER* - invalid parameter,
- *NMDL\_RESULT\_NO\_BOARD* – no board,
- *NMDL\_RESULT\_BOARD\_RESET\_ERROR* – board reset error,
- *NMDL\_RESULT\_INIT\_CODE\_LOADING\_ERROR* – initialization code loading error,
- *NMDL\_RESULT\_CORE\_HANDLE\_RETRIEVAL\_ERROR* – error of obtaining the identifier of the board,
- *NMDL\_RESULT\_FILE\_LOADING\_ERROR* – program loading error,
- *NMDL\_RESULT\_MEMORY\_WRITE\_ERROR* – write memory error,
- *NMDL\_RESULT\_MEMORY\_READ\_ERROR* – read memory error,
- *NMDL\_RESULT\_MEMORY\_ALLOCATION\_ERROR* – memory allocation error,
- *NMDL\_RESULT\_MODEL\_LOADING\_ERROR* – neural network model loading error,

- *NMDL\_RESULT\_INVALID\_MODEL* – wrong neural network model,
- *NMDL\_RESULT\_BOARD\_SYNC\_ERROR* – synchronization error,
- *NMDL\_RESULT\_BOARD\_MEMORY\_ALLOCATION\_ERROR* – allocation memory error on the board,
- *NMDL\_RESULT\_NN\_CREATION\_ERROR* – model creation error on the board,
- *NMDL\_RESULT\_NN\_LOADING\_ERROR* – neural network model loading error,
- *NMDL\_RESULT\_NN\_INFO\_RETRIEVAL\_ERROR* – neural network model information obtaining error,
- *NMDL\_RESULT\_MODEL\_IS\_TOO\_BIG* – neural network model is too big,
- *NMDL\_RESULT\_NOT\_INITIALIZED* – NMDL is not initialized,
- *NMDL\_RESULT\_BUSY* – board is busy (frame is processing),
- *NMDL\_RESULT\_UNKNOWN\_ERROR* – unknown error.

## 9.1.5. NMDL\_Tensor

Structure describing a tensor.

```
typedef struct tagNMDL_Tensor {
    unsigned int width;
    unsigned int height;
    unsigned int depth;
} NMDL_Tensor;
```

- *width* - tensor width,
- *height* - tensor height,
- *depth* – tensor depth.

## 9.2. Functions

### 9.2.1. NMDL\_Blink

LED indication for module identification.

```
NMDL_RESULT NMDL_Blink(
    unsigned int board_type,
    unsigned int board_number
);
```

- *board\_type* - [in] type of module on which the LED indication procedure is called. One of the enumeration values [NMDL\\_BOARD\\_TYPE](#).

- *board\_number* - [in] the serial number of the module on which the LED indication procedure is called.

## 9.2.2. NMDL\_Create

Instantiate the NMDL and get the *NMDL* instance id.

```
NMDL_RESULT NMDL_Create(
    NMDL_HANDLE *nmdl
);
```

- *nmdl* - [out] instance id *NMDL*.

After working with the NMDL instance, you need to release the allocated resources by calling [NMDL\\_Destroy](#).

## 9.2.3. NMDL\_Destroy

Deleting an NMDL Instance.

```
void NMDL_Destroy(
    NMDL_HANDLE nmdl
);
```

- *nmdl* - [in] instance id *NMDL*.

The function is called to release resources allocated by [NMDL\\_Create](#) and [NMDL\\_Initialize](#) calls.

## 9.2.4. NMDL\_GetBoardCount

Request for the number of detected modules of a given type.

```
NMDL_RESULT NMDL_GetBoardCount(
    unsigned int board_type,
    unsigned int *boards
);
```

- *board\_type* - [in] the type of modules being polled. One of the enumeration values [NMDL\\_BOARD\\_TYPE](#).
- *boards* - [out] the number of modules found.

For MC127.05 simulator (type *NMDL\_BOARD\_TYPE\_SIMULATOR*) one module is always detected.

## 9.2.5. NMDL\_GetLibVersion

Library version query *NMDL*.

```
NMDL_RESULT NMDL_GetLibVersion(
```

```

    unsigned int *major,
    unsigned int *minor,
    unsigned int *patch
);

```

- *major* - [out] major version number.
- *minor* - [out] minor version number.
- *patch* - [out] patch version number.

## 9.2.6. NMDL\_GetModelInfo

Request information about the model.

```

NMDL_RESULT NMDL_GetModelInfo(
    NMDL_HANDLE nmdl,
    NMDL_ModelInfo *model_info
);

```

- *nmdl* - [in] descriptor *NMDL*
- *model\_info* - [out] model information.

## 9.2.7. NMDL\_GetOutput

Request for the processing result.

```

NMDL_RESULT NMDL_GetOutput(
    NMDL_HANDLE nmdl,
    unsigned int batch_num,
    float *output,
    double *fps
);

```

- *nmdl* - [in] instance id *NMDL*.
- *batch\_num* - [in] cluster number for which the processing result is requested. Relevant only when working with *MC127.05* modules, *NMCard* or a simulator. When working with *MC121.01* and *NMStick*, set 0.
- *output* - [out] pointer to the output tensor buffer. The output tensors are placed one after another in the order of their appearance in the processing graph, that is, ordered by levels in the graph. Tensor parameters are defined in the [NMDL\\_ModelInfo](#) structure. The size of the required buffer is returned in the [NMDL\\_RequiredOutputFloats](#) function. May be 0.
- *fps* - [out] processing performance (frames per second). The performance value is defined only for cluster 0 (*batch\_num* = 0), for other clusters it writes 0.0f. May be 0.

## 9.2.8. NMDL\_GetStatus

Request processing status. The function is called to check that the frame has been processed.

```
NMDL_RESULT NMDL_GetStatus(
    NMDL_HANDLE nmdl,
    unsigned int batch_num,
    unsigned int *status
);
```

- *nmdl* - [in] instance id *NMDL*
- *batch\_num* - [in] cluster number for which the processing status is requested. Relevant only when working with *MC127.05* modules, *NMCard* or a simulator. When working with *MC121.01* and *NMStick*, set 0.
- *status* - [out] the state of the cluster at the time the function is called. One of the enumeration values [NMDL\\_PROCESS\\_FRAME\\_STATUS](#).

## 9.2.9. NMDL\_Initialize

NMDL initialization.

```
NMDL_RESULT NMDL_Initialize(
    NMDL_HANDLE nmdl,
    unsigned int board_type,
    unsigned int board_number,
    const float *model,
    unsigned int model_floats,
    unsigned int use_batch_mode
);
```

- *nmdl* - [in] instance id *NMDL*.
- *board\_type* - [in] the type of the module being initialized. One of the enumeration values [NMDL\\_BOARD\\_TYPE](#).
- *board\_number* - [in] the sequence number of the module being initialized.
- *model* - [in] pointer to the buffer containing the model in the format *nm7* (for *MC121.01* and *NMStick*) or *nm8* (for *MC127.05*, *NMCard* and simulator).
- *model\_floats* - [in] the size of the model in real numbers with single precision.
- *use\_batch\_mode* - [in] flag for using batch processing mode. Relevant only for *MC127.05* modules, *NMCard* or a simulator. When working with *MC121.01* and *NMStick*, set 0.

After working with NMDL functions, you need to release the allocated resources by calling the deinitialization function [NMDL\\_Release](#).

## 9.2.10. NMDL\_ProcessFrame

Single frame processing.

```
NMDL_RESULT NMDL_ProcessFrame(
    NMDL_HANDLE nmdl,
    unsigned int batch_num,
    const float *frame
```

```
) ;
```

- *nmdl* - [in] instance id *NMDL*
- *batch\_num* - [in] cluster number on which processing is started. Relevant only when working with *MC127.05* modules, *NMCard* or a simulator. When working with *MC121.01* and *NMStick*, set 0.
- *frame* - [in] pointer to the buffer of the processed frame.

### 9.2.11. NMDL\_Release

Deinitialize NMDL.

```
void NMDL_Release(
    NMDL_HANDLE nmdl
);
```

- *nmdl* - [in] instance id *NMDL*

### 9.2.12. NMDL\_RequiredOutputFloats

Returns the size of the buffer in *float32* elements required to accommodate the result (*output*) when calling [NMDL\\_GetOutput](#).

```
unsigned int NMDL_RequiredOutputFloats(
    const NMDL_ModelInfo *model_info
);
```

- *model\_info* - [out] information about the model previously obtained by calling [NMDL\\_GetModelInfo](#).

# 10. Description of identifiers and functions nmdl\_compiler

## 10.1. Identifiers

### 10.1.1. NMDL\_COMPILER\_BOARD\_TYPE

Module types.

```
typedef enum tagNMDL_COMPILER_BOARD_TYPE {
    NMDL_COMPILER_BOARD_TYPE_MC12101,
    NMDL_COMPILER_BOARD_TYPE_MC12705
} NMDL_COMPILER_BOARD_TYPE;
```

- *NMDL\_COMPILER\_BOARD\_TYPE\_MC12101* - modules *MC121.01* и *NMStick*.
- *NMDL\_COMPILER\_BOARD\_TYPE\_MC12705* - modules *MC127.05*, *NMCard* and simulator.

### 10.1.2. NMDL\_COMPILER\_RESULT

Returned result.

```
typedef enum tagNMDL_COMPILER_RESULT {
    NMDL_COMPILER_RESULT_OK,
    NMDL_COMPILER_RESULT_MEMORY_ALLOCATION_ERROR,
    NMDL_COMPILER_RESULT_MODEL_LOADING_ERROR,
    NMDL_COMPILER_RESULT_INVALID_PARAMETER,
    NMDL_COMPILER_RESULT_INVALID_MODEL,
    NMDL_COMPILER_RESULT_UNSUPPORTED_OPERATION
} NMDL_COMPILER_RESULT;
```

- *NMDL\_COMPILER\_RESULT\_OK* - no errors,
- *NMDL\_COMPILER\_RESULT\_MEMORY\_ALLOCATION\_ERROR* – memory allocation error,
- *NMDL\_COMPILER\_RESULT\_MODEL\_LOADING\_ERROR* – error loading neural network model,
- *NMDL\_COMPILER\_RESULT\_INVALID\_PARAMETER* - invalid parameter,
- *NMDL\_COMPILER\_RESULT\_INVALID\_MODEL* – error in the model,
- *NMDL\_COMPILER\_RESULT\_UNSUPPORTED\_OPERATION* – the model contains an unsupported operation.

## 10.2. Functions

### 10.2.1. NMDL\_COMPILER\_CompileDarkNet

Compilation of the original model in DarkNet format.

```
NMDL_COMPILER_RESULT NMDL_COMPILER_CompilerDarkNet(
    unsigned int board,
    const char* src_model,
    unsigned int src_model_size,
    const char* src_weights,
    unsigned int src_weights_size,
    float** dst_model,
    unsigned int* dst_model_floats
);
```

- *board* - [in] module type identifier for which the model is compiled. One of the enumeration values [NMDL\\_COMPILER\\_BOARD\\_TYPE](#).
- *src\_model* - [in] buffer of the original model in the *DarkNet* format, previously read from the *.cfg* file.
- *src\_model\_size* - [in] the size of the original model buffer in bytes.
- *src\_weights* - [in] coefficient buffer previously read from *.weights* .
- *src\_weights\_size* - [in] size of the coefficient buffer in bytes.
- *dst\_model* - [out] compiled model buffer. Memory allocation occurs in a function.
- *dst\_model\_floats* - [out] compiled model buffer size in *float32* .

### 10.2.2. NMDL\_COMPILER\_CompilerONNX

Compilation of the source model in ONNX format.

```
NMDL_COMPILER_RESULT NMDL_COMPILER_CompilerONNX(
    unsigned int board,
    const char* src_model,
    unsigned int src_model_size,
    float** dst_model,
    unsigned int* dst_model_floats
);
```

- *board* - [in] module type identifier for which the model is compiled. One of the enumeration values [NMDL\\_COMPILER\\_BOARD\\_TYPE](#).
- *src\_model* - [in] the buffer of the original model in the *ONNX* format, previously read from the file *.pb* .
- *src\_model\_size* - [in] the size of the original model buffer in bytes.
- *dst\_model* - [out] compiled model buffer. Memory allocation occurs in a function.

- *dst\_model\_floats* - [out] compiled model buffer size in *float32* .

### **10.2.3. NMDL\_COMPILER\_FreeModel**

Freeing compile-time allocated memory.

```
NMDL_COMPILER_RESULT NMDL_COMPILER_FreeModel (
    unsinged int board,
    char* dst_model
);
```

- *board* - [in] module type identifier for which the model was compiled. One of the enumeration values [NMDL\\_COMPILER\\_BOARD\\_TYPE](#).
- *dst\_model* - [in]freed memory buffer. Buffer memory was allocated by calling [NMDL\\_COMPILER\\_CompileDarkNet](#) or [NMDL\\_COMPILER\\_CompileONNX](#).

### **10.2.4. NMDL\_COMPILER\_GetLastError**

Returns a constant string describing the last error.

```
const char *NMDL_COMPILER_GetLastError();
```

# 11. Description of identifiers and functions nmdl\_image\_converter

## 11.1. Identifiers and structures

### 11.1.1. NMDL\_IMAGE\_CONVERTER\_BOARD\_TYPE

Types of modules.

```
typedef enum tagNMDL_IMAGE_CONVERTER_BOARD_TYPE {
    NMDL_IMAGE_CONVERTER_BOARD_TYPE_MC12101,
    NMDL_IMAGE_CONVERTER_BOARD_TYPE_MC12705
} NMDL_IMAGE_CONVERTER_BOARD_TYPE;
```

- *NMDL\_IMAGE\_CONVERTER\_BOARD\_TYPE\_MC12101* - modules *MC121.01* и *NMStick*.
- *NMDL\_IMAGE\_CONVERTER\_BOARD\_TYPE\_MC12705* - modules *MC127.05*, *NMCard* and simulator.

### 11.1.2. NMDL\_IMAGE\_CONVERTER\_COLOR\_FORMAT

Pixel format identifiers. Describes the ordering of RGB color components in a pixel.

```
typedef enum tagNMDL_IMAGE_CONVERTER_COLOR_FORMAT {
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_RGB,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_RGB,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_GRB,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_GBR,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_BRG,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_BGR,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_INTENSITY,
} NMDL_IMAGE_CONVERTER_COLOR_FORMAT;
```

## 11.2. Functions

### 11.2.1. NMDL\_IMAGE\_CONVERTER\_Convert

Image preparation.

```
int NMDL_IMAGE_CONVERTER_Convert(
    const char* src,
    float* dst,
    unsigned int src_size,
    unsigned int dst_width,
    unsigned int dst_height,
    unsigned int dst_color_format,
    const float rgb_divider[3],
    const float rgb_adder[3],
    unsigned int board_type
);
```

- *src* - [in] buffer with the original image. The contents of the buffer correspond to the contents of the image file. Images can be in .bmp, .gif, .jpg and .png formats.
- *dst* - [out] buffer for the prepared image. The buffer memory must be pre-allocated. The buffer size must not be less than the value returned by the [NMDL\\_IMAGE\\_CONVERTER\\_RequiredSize function](#).
- *src\_size* - [in] the size of the source image buffer in bytes.
- *dst\_width* - [in] width of the prepared image.
- *dst\_height* - [in] the height of the prepared image.
- *dst\_color\_format* - [in] pixel format identifier of the prepared image. One of the enumeration values [NMDL\\_IMAGE\\_CONVERTER\\_COLOR\\_FORMAT](#).
- *rgb\_divider* - [in] divider in expression  $dst = src / divider + adder$ . The above operation is performed on each channel of each pixel of the original image. *rgb\_divider[0]* - red channel divider, *rgb\_divider[1]* - green channel divider, *rgb\_divider[2]* - blue channel divider.
- *rgb\_adder* - [in] term in expression  $dst = src / divider + adder$ . The above operation is performed on channels of each pixel of the original image. *rgb\_adder[0]* - red channel adder, *rgb\_adder[1]* - green channel adder, *rgb\_adder[2]* - blue channel adder.
- *board\_type* - [in] the identifier of the type of computing module on which the processing is supposed. One of the enumeration values [NMDL\\_IMAGE\\_CONVERTER\\_BOARD\\_TYPE](#).

For grayscale images, when *dst\_color\_format* = [NMDL\\_IMAGE\\_CONVERTER\\_COLOR\\_FORMAT\\_INTENSITY](#), only the divisor *rgb\_divider[0]* and the *rgb\_adder[0]* term are used. Other divisors (*rgb\_divider[1]* and *rgb\_divider[2]*) and terms (*rgb\_adder[1]* and *rgb\_adder[2]*) are ignored and can be set to any value.

Returned value: 0 - normal completion, -1 - error.

## 11.2.2. NMDL\_IMAGE\_CONVERTER\_RequiredSize

Returns the size of the buffer in *float32* elements to hold the prepared image.

```
int NMDL_IMAGE_CONVERTER_RequiredSize(
    unsigned int dst_width,
    unsigned int dst_height,
    unsigned int dst_color_format,
    unsigned int board_type
);
```

- *dst\_width* - [in] width of the prepared image.
- *dst\_height* - [in] the height of the prepared image.
- *dst\_color\_format* - [in] pixel format identifier of the prepared image. One of the enumeration values [NMDL\\_IMAGE\\_CONVERTER\\_COLOR\\_FORMAT](#).

- *board\_type* - [in] the identifier of the type of computing module on which the processing is supposed. One of the enumeration values [NMDL\\_IMAGE\\_CONVERTER\\_BOARD\\_TYPE](#).

Returned value: 0 - normal completion, -1 - error.



**Research Centre Module**  
**Box: 166, Moscow, 125190, Russia**  
**Tel: +7 (499) 152-9698**  
**Fax: +7 (499) 152-4661**  
**E-Mail: sales@module.ru**  
**WWW: <http://www.module.ru>**

©RC "Module", 2021

All rights reserved.

Neither the whole nor any part of the information contained in, or the product described in this overview may be adapted or reproduced in any form except with the prior written permission of the copyright holder.

RC Module reserves the right to make changes without further notices to product herein to improve reliability, function or design. RC Module shall not be liable for any loss or damage arising from the use of any information in this overview or any error or omission in such information, or any incorrect use of the product.

Printed in Russia Data of issue: