

NMDL - Программное обеспечение для реализации глубоких нейронных сетей на платформе NeuroMatrix®

NMDL Руководство пользователя



Содержание

1. Общие сведения	. 6
1.1. Назначение	. 6
1.2. Быстрый старт	. 6
1.3. Производительность NMDL	6
2. Состав ПО	. 9
3. Установка	11
3.1. Установка в Windows	11
3.2. Установка в Linux	11
3.3. Дополнительные компоненты	11
3.4. Подготовка модуля МС121.01	12
3.5. Подготовка устройства NMStick	13
3.6. Подготовка модулей MC127.05 и NMCard	14
4. Компиляция модели	15
4.1. Поддерживаемые операции	16
5. Подготовка изображений	18
6. Режимы обработки	20
7. Демонстрационная программа	23
8. Пример использования NMDL	30
9. Описание идентификаторов, функций и структур NMDL	36
9.1. Идентификаторы и структуры	36
9.1.1. NMDL BOARD TYPE	36
9.1.2. NMDL_ModelInfo	36
9.1.3. NMDL_PROCESS_FRAME_STATUS	36
9.1.4. NMDL_RESULT	37
9.1.5. NMDL_Tensor	38
9.2. Функции	38
9.2.1. NMDL Blink	38
9.2.2. NMDL Create	39
9.2.3. NMDL Destroy	39
9.2.4. NMDL GetBoardCount	39
9.2.5. NMDL GetLibVersion	40
9.2.6. NMDL GetModelInfo	40
9.2.7. NMDL GetOutput	40
9.2.8. NMDL GetStatus	41
9.2.9. NMDL Initialize	41
9.2.10. NMDL ProcessFrame	42
9.2.11. NMDL Release	42
9.2.12. NMDL RequiredOutputFloats	42
10. Описание идентификаторов и функций nmdl compiler	43
10.1. Идентификаторы	43
10.1.1. NMDL COMPILER BOARD TYPE	43
10.1.2. NMDL_COMPILER_RESULT	43
10.2. Функции	44
10.2.1. NMDL_COMPILER_CompileDarkNet	44
10.2.2. NMDL_COMPILER_CompileONNX	44



10.2.3. NMDL COMPILER FreeModel	45
10.2.4. NMDL_COMPILER_GetLastError	45
11. Описание идентификаторов и функций nmdl_image_converter	46
11.1. Идентификаторы и структуры	46
11.1.1. NMDL_IMAGE_CONVERTER_BOARD_TYPE	46
11.1.2. NMDL_IMAGE_CONVERTER_COLOR_FORMAT	46
11.2. Функции	46
11.2.1. NMDL_IMAGE_CONVERTER_Convert	46
11.2.2. NMDL_IMAGE_CONVERTER_RequiredSize	47



Список иллюстраций

3.1. Перемычки на модуле МС121.01	13
6.1. Кластеры СБИС К1879ВМ8Я.	20
6.2. Режимы обработки	21
7.1. Внешний вид программы в процессе работы	24
7.2. Внешний вид окна выбора модуля	25
7.3. Внешний вид окна результатов классификации	28
7.4. Отображение результатов детектирования	29

Список таблиц

1.1. FPS	 !
1.2. Latency (ms)	 ;



1. Общие сведения

1.1. Назначение

Программный модуль *NMDL* позволяет запускать предварительно обученную глубокую сверточную нейронную сеть на вычислительных модулях *MC121.01*, *MC127.05*, *NMStick*, *NMCard* и на симуляторе модуля *MC127.05*. Программный модуль состоит из 2 частей. Одна часть работает на персональном компьютере (хост) под управлением 64 разрядных OC Microsoft® Windows 7/10 или Linux, другая часть запускается и работает на процессоре вычислительного модуля. Связь устройств *MC121.01* и *NMStick* с хостом осуществляется по каналу USB2.0, для связи модулей *MC127.05* и *NMCard* с хостом используется интерфейс PCIe.

1.2. Быстрый старт

В разделе приводится пошаговая инструкция для быстрого начала работы с <u>демонстрационной программой</u>. Более полную информацию о *NMDL* можно получить в соответсвующих разделах руководства. Здесь описывается демонстрация обработки изображения с применением нейронной сети для классификации *Squeezenet* на симуляторе устройства на базе СБИС 1879ВМ8Я.

- Установите дистрибутив *NMDL* так, как описано в разделе <u>Установка</u>.
- Перейдите в каталог *bin* в директории установки *NMDL* (по-умолчанию в Windows "c:\Program Files\Module\NMDL", по-умолчанию в Linux "/opt/nmdl") и запустите демонстрационную программу *nmdl_gui*.
- Выберите устройство с помощью меню *File Open Board*.... Убедитесь, что в диалоговом окне выбран симулятор.
- Выберите описание модели с помощью меню *File Open Description....* В диалоговом окне выбора файла выберите файл *PATH_TO_NMDL/nmdl_ref_data/* squeezenet_imagenet/description08.xml.
- Выберите обрабатываемое изображение с помощью меню *File Open Picture....* В диалоговом окне выбора файла выберите файл *PATH_TO_NMDL/nmdl_ref_data/ squeezenet_imagenet/frame.bmp.*
- Запустите обработку с помощью меню *File Run*. В результате обработки всплывёт окно классификации с вычисленными вероятностями в порядке уменьшения. Верный класс для выбранного изображения "lakeside".

1.3. Производительность NMDL

В таблицах приводятся значения производительности (кадры в секунду FPS) и значения задержки от начала обработки кадра до получения результата (Latency). Рядом с названием сети в скобках указан размер обрабатываемого изображения.



Таблица 1.1 - FPS

	MC121.01	NMStick	MC127.05 и NMCard	MC127.05 и NMCard batch- mode*
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
squeezenet (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

	MC121.01	NMStick	MC127.05 NMCard	и MC127.05 и NMCard batch- mode*
alexnet (227x227)	290	302	79	308
inception v2 (299x299)	3 1587	1653	78	197
inception v3 (512x512)	3 4166	4340	254	735
resnet 18 (224x224)	3 439	457	40	85
resnet 50 (224x224)) 1250	1300	82	194
squeezenet (224x224)	120	125	13	40
u-net (512x512)**	-	-	500	2000
yolo v2 tiny (416x416)	/ 862	898	47	132
yolo v2 (416x416)	3 10000	10416	270	889
yolo v3 tiny (416x416)	694	725	36	113

Таблица 1.2 - Latency (ms)

* Описание *batch-mode* приводится в разделе <u>"Режимы обработки"</u>

** В модели u-net произведена замена слоёв transposed_convolution на upsampling.

2. Состав ПО

ПО реализации нейронных сетей состоит из программных модулей (API), утилит и руководства.

Файлы АРІ для разработки программ с использованием *NMDL*:

- *nmdl.dll/nmdl.so* программный модуль для применения обученной нейронной сети. См. раздел <u>"Описание идентификаторов, функций и структур nmdl"</u>.
- *nmdl.lib* библиотека для раннего связывания программ с *NMDL* в среде MSVC++.
- *nmdl.h* заголовочный файл с описанием структур и функций АРІ.
- *nmdl_compiler.dll/nmdl_compiler.so* программный модуль компилятор моделей *ONNX/DarkNet* во внутреннее представление. См. <u>"Описание идентификаторов и</u> <u>функций nmdl_compiler"</u>
- *nmdl_compiler.lib* библиотека для раннего связывания модуля компилятора моделей в среде MSVC++.
- *nmdl_compiler.h* заголовочный файл с описанием структур и функций компилятора моделей.
- *nmdl_image_converter.dll/nmdl_image_converter.so* программный модуль для подготовки обрабатываемых изображений. См. <u>"Описание идентификаторов и функций nmdl_image_converter"</u>
- *nmdl_image_converter.dll / nmdl__image_converter.so* a program module for preparing processed images. See <u>"Description of nmdl_image_converter identifiers and functions"</u>
- *nmdl_image_converter.lib* модуль для раннего связывания модуля подготовки изображений в среде MSVC++.
- *nmdl_image_converter.h* заголовочный файл с описанием структур и функций для подготовки изображений.

Заголовочные файлы и библиотеки раннего связывания размещаются в каталогах *include* и *lib* директории *NMDL*.

Утилиты:

• *nmdl_compiler_console* - утилита командной строки для компиляции моделей из форматов ONNX и DarkNet во внутренний формат для загрузки на вычислительные модули. Файл модели ONNX обычно имеет расширение .pb. Модель в формате DarkNet сохраняется в двух файлах - с расширением .cfg и расширением .weights. Подготовленная модель для устройств MC121.01 и NMStick имеет расширение .nm7. Модель для MC127.05 и NMCard имеет расширение .nm8. См. раздел "Компиляция модели".



- *nmdl_nmdl_image_converter_console* утилита командной строки для подготовки обрабатываемых изображений. См. раздел <u>"Подготовка изображений"</u>.
- *nmdl_gui* оконная утилита для демонстрации функциональных возможностей NMDL. См. раздел <u>"Демонстрационная программа"</u>.

3. Установка

3.1. Установка в Windows

NMDL распространяется в виде установочного дистрибутива. Поддерживаются только 64-х битовые системы.

Для установки дистрибутива запустите исполняемый файл инсталлятора с правами администратора. Следуйте инструкциям мастера установкии.

Для работы с программными модулями необходимо включить их в область "видимости" операционной системы. Одно из решений - создание переменной среды окружения, например, с именем *NMDL*, в которой записывается путь к каталогу с заголовочными и бинарными файлами, и добавление созданной переменной к переменной окружения РАТН.

Для использования модуля nmdl в C/C++ программах включите в исходный файл директиву #include "nmdl.h" и свяжите программу с библиотекой NMDL. Если используется среда разработки MSVC++, то для раннего связывания подключите к проекту файл nmdl.lib. Можно создать и использовать переменную окружения NMDL для задания пути к файлам nmdl.h и nmdl.lib. Таким же образом можно подключить модули nmdl_compiler и nmdl_image_converter.

3.2. Установка в Linux

Распространяется в виде .deb пакета. Поддерживаются только 64-х битовые системы семейства Debian.

Для установки используйте менеджер пакетов dpkg.

Например:

dpkg -i NMDL.deb

3.3. Дополнительные компоненты

Вместе с программным модулем можно получить дополнительные компоненты для проверки и демонстрации работы *NMDL*.

Компоненты распространяются в архивах:

- *nmdl_ref_data_alexnet_imagenet.zip* архив для демонстрации работы сети ALEXNET.
- *nmdl_ref_data_inception_v3_imagenet.zip* архив для демонстрации работы сети *INCEPTION V3*.
- *nmdl_ref_data_resnet_18_imagenet.zip* архив для демонстрации работы сети *RESNET18*.



- *nmdl_ref_data_resnet_50_imagenet.zip* архив для демонстрации работы сети *RESNET50*.
- *nmdl_ref_data_squeezenet_imagenet.zip* архив для демонстрации работы сети *SQUEEZENET*.
- *nmdl_ref_data_unet_no_transp_conv_covid.zip* архив для демонстрации работы сети *U-NET*. В модели произведена замена слоёв *transposed_convolution* на *upsampling*.
- *nmdl_ref_yolo2_tiny_pascal_voc.zip* архив для демонстрации работы сети YOLONET V2 TINY.
- *nmdl_ref_yolo3_coco.zip* архив для демонстрации работы сети YOLONET V3.
- *nmdl_ref_yolo3_tiny_coco.zip* архив для демонстрации работы сети YOLONET V3 *TINY*.

В каталогах содержатся файлы:

- *frame.bmp* тестовое изображение.
- *model.cfg* модель в формате *DARKNET*.
- *model.pb* модель в формате ONNX.
- model.weights весовые коэффициенты модели в формате DARKNET.
- *description07.xml* файл описания модели для устройств *MC121.01* и *NMStick* для запуска в программе nmdl_gui.
- *description08.xml* файл описания модели для модулей *MC127.05*, *NMCard* и симулятора для запуска в программе nmdl_gui.

Для подготовки демонстрационных данных, которые будут обрабатываться на вычислительном модуле необходимо распаковать их в каталог *nmdl_ref_data*, выполнить компиляцию моделей и подготовку изображений так, как это описывается в разделах <u>"Компиляция модели"</u> и <u>"Подготовка изображений"</u> данного руководства. Для удобства компиляции в архив включены скрипты *prepare.cmd* и *prepare.sh*. В результате работы скрипта в каждом каталоге появятся файлы:

- *frame07* подготовленное для обработки на *MC121.01* и *NMStick* изображение.
- *frame08* подготовленное для обработки на MC127.05 и NMCard изображение.
- model.nm7 компилированная модель для загрузки на MC121.01 и NMStick.
- *model.nm8* компилированная модель для загрузки на *MC127.05* и *NMCard*.

3.4. Подготовка модуля МС121.01

При использовании NMDL с устройствами *MC121.01* необходимо до подключения модуля к USB-разъёму ПК установить следующие перемычки (см. рисунок <u>3.1</u>):



- Контакт 1-2 разъема Х9.
- Контакт 3-4 разъема Х9.
- Контакт 5-6 разъема Х9.
- При использовании питания от USB контакт 1-2 разъема X11 (при использовании блока питания контакт разомкнуть).
- Контакт 1-2 разъема Х18.



Рисунок 3.1 - Перемычки на модуле МС121.01

3.5. Подготовка устройства NMStick

При использовании NMDL с *NMStick* необходимо выполнить инсталляцию ПО поддержки устройства (входит в комплект поставки изделия) и подключить устройство к USB-разъёму ПК.



3.6. Подготовка модулей MC127.05 и NMCard

Для работы с NMDL необходимо установить плату в свободный слот PCIe и выполнить инсталляцию ПО поддержки модуля, входящее в комплект поставки изделия.



4. Компиляция модели

Нейронная сеть должна быть предварительно обучена с помощью нейросетевого пакета (например, Microsoft® Cognitive Toolkit (CNTK)), и преобразована к форматам *ONNX* или *DarkNet*.

Исходная модель в формате ONNX, как правило, хранится в файле со структурой Google Protocol Buffer и имеет расширение *.pb или *.onnx.

Модели для сетей типа YOLO могут храниться в формате DarkNet. В этом случае модель описывается двумя файлами - файл с расширением *.cfg содержит описание графа обработки, файл *.weights содержит веса для свёрток.

Файлы моделей ONNX и DarkNet являются результатом обучения нейронной сети и одновременно исходными данными для компилятора, который в результате обработки создаёт файлы с расширением *nm7* для устройств *MC121.01* и *NMStick*, или *nm8* для модулей *MC127.05*, *NMCard* и симулятора.

Модели *nm*7 и *nm*8 являются бинарными образами скомпилированных пользовательских моделей. Их структура не документируется и зависит от целевого устройства и версии компилятора.

Компилятор выполнен в виде динамически загружаемого модуля и может быть встроен в программу пользователя. Можно также воспользоваться консольной утилитой nmdl_compiler_console для выполнения преобразования моделей из командной строки.

nmdl_compiler_console BOARD_TYPE NN_TYPE
 SRC_FILENAME DST_FILENAME [WEIGHTS_FILENAME]

Аргументы командной строки:

- *BOARD_TYPE* тип модуля. Допустимые значения: "*MC12101*" преобразование модели в формат nm7 для устройств *MC121.01* и *NMStick*, "*MC12705*" преобразование модели в формат nm8 для модулей *MC127.05*, *NMCard* и симулятора.
- *NN_TYPE* формат файла входной модели. Допустимые значения: "ONNX" или "DARKNET".
- *SRC_FILENAME* имя файла исходной модели.
- *DST_FILENAME* имя файла выходной модели.
- *WEIGHTS_FILENAME* имя файла с весовыми коэффициентами модели. Нужен только для моделей в формате DarkNet.

Пример:



>nmdl_compiler_console MC12705 ONNX
 squeezenet.pb squeezenet.nm8
>nmdl_compiler_console MC12705 DARKNET

4.1. Поддерживаемые операции

Компилятором поддерживается следующий набор операций:

- 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 (при условии, что операция выполняется сразу после свёртки).
- 5. Сопсат (по осям каналов и ширины, количество входных тензоров не более 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 может работать с моделями с одним обрабатываемым изображением, то есть обрабатывается один входной тензор.

Количество терминальных узлов (выходные тензоры) - не более 5.

5. Подготовка изображений

Обрабатываемые изображения необходимо предварительно сконвертировать в вещественный формат (тензор). Сделать это можно с помощью утилиты *nmdl_image_converter_console*

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

Аргументы командной строки:

- SRC_FILE имя входного файла (bmp, gif, jpg, emf, png, tiff),
- DST FILE имя выходного файла,
- *W* ширина тензора изображения,
- *Н* высота тензора изображения,
- *F* порядок следования цветовых каналов в тензоре изображения (допустимые значения: *rgb*, *rbg*, *grb*, *gbr*, *brg*, *bgr*, intensity один канал градации серого),
- DR делитель для красного канала пикселя в выражении dst = src / D + A (вещественная величина),
- DG делитель для зелёного канала пикселя в выражении dst = src / D + A (вещественная величина),
- DB делитель для голубого канала пикселя в выражении dst = src / D + A (вещественная величина),
- AR слагаемое для красного канала пикселя в выражении dst = src / D + A (вещественная величина).
- AG слагаемое для зелёного канала пикселя в выражении dst = src / D + A (вещественная величина).
- AB слагаемое для голубого канала пикселя в выражении dst = src / D + A (вещественная величина).
- *BT* тип модуля на котором предполагается обработка (допустимые значения: mc12101, mc12705).

Для изображений с градацией серого, когда аргумент F задан intensity, используются только делитель DR и слагаемое AR. Остальные делители (DG и DB) и слагаемые (AG и AB) игнорируются и могут быть установлены в любые значения.

Программа масштабирует входное изображение относительно центра в соответствии с заданными аргументами.

Подготовленные изображения имеют либо планарный (для *MC121.01* и *NMStick*), либо пиксельный (для *MC127.05* и *NMCard*) формат расположения элементов типа *float32*.



В планарном формате в файл записываются поочерёдно плоскости каждого из каналов. В этом случае буфер можно представить как массив в стиле C/C++:

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

То есть самый быстро меняющийся индекс - это индекс по ширине изображения. Количество каналов - три (RGB каналы), или один для одноканальных изображений (градации серого). В подготовленном буфере размер выравнивается по ширине изображения. Для изображений с чётной шириной размер буфера составит: *size* = *width* * *height* * *channels* (float32), для изображений с нечётной шириной: *size* = (*width* + 1) * *height* * *channels* (float32).

В пиксельном формате в файл записываются последовательно каналы пикселей изображения. Буфер можно представить как массив в стиле C/C++:

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

То есть самый быстро меняющийся индех - это индекс по каналам. Количество каналов - три (RGB каналы), или один для одноканальных изображений (градации серого). В подготовленном буфере размер выравнивается по каналам изображения до ближайшего чётного значения. Размер буфера с подготовленным изображением для модулей MC127.05 и NMCard составит: size = width * height * (channels + 1) (float32).





6. Режимы обработки

Архитектурные особенности используемых в вычислительных модулях СБИС позволяют реализовать нейросетевые алгоритмы по разному. Разнообразие определяется аппаратными средствами для выполнения параллельных вычислений. В реализации *NMDL* использовалась глубокая оптимизация алгоритмов и выбирались наиболее быстрые способы параллелизма, инкапсулируюя детали разработки. Тем не менее, в реализации *NMDL* для модуля MC127.05 остаётся вариативность решений, которую целесообразно вынести в интерфейс и предоставить пользователю возможность выбрать тот, или иной способ обработки в зависимости от решаемой задачи.

Используемая в модулях *MC127.05* и *NMCard* СБИС К1879ВМ8Я имеет кластерную организацию.



Рисунок 6.1 - Кластеры СБИС К1879ВМ8Я.

На рисунке <u>6.1</u>:



- *PC* процессорные кластеры;
- *NMPU* процессорные узлы NMC4;
- *CPU* кластерные процессорные ядра ARM Cortex-A5;
- *ССРU* центральное процессорное ядро ARM Cortex-A5.

Внутри кластеров объём вычислений и данных распределяются на четыре процессорных узла NMC4. Здесь используются тесные аппаратные связи при межпросессорном взаимодействии внутри кластеров и выполнена глубокая оптимизация над функциями-примитивами.

Обработка на кластерах же может быть организована двумя способами: с разделением данных между кластерами и работа в пакетном режиме (*batch-mode*).



Рисунок 6.2 - Режимы обработки

При обработке с разделением данных (рисунок <u>6.2</u> A)) можно добиться минимальной латентности - времени от начала обработки до момента получения результата. Однако в этом случае неизбежны накладные расходы по организации параллелизма. Например, при выполнении свёрток над разделёнными тензорами необходимо компенсировать обработку границ, выполнять переупаковки, транзит данных между кластерами и т. п. Производительность - кадры в секунду - *FPS* = 1/Latency.



При пакетной обработке (рисунок <u>6.2</u> В)) на каждом кластере обрабатывается один кадр. Кластеры выполняют одинаковую и независимую обработку. В этом случае накладные расходы по организации распределения данных отсутствует и достигается максимальная производительность, при этом увеличивается латентность. Производительность - кадры в секунду - FPS = 4/Latency.

Управление режимами обработки возможно только при работе с модулями *MC127.05*, *NMCard* или с симулятором.



7. Демонстрационная программа

Для демонстрации функциональных возможностей библиотеки NMDL была разработана утилита nmdl_gui. В задачи утилиты входят:

- Инициализация библиотеки NMDL;
- Обнаружение и идентификация доступных ускорителей (*симулятор*, MC121.01, MC127.05, NMStick, NMCard);
- Загрузка файла модели нейронной сети (.*nm7* или .*nm8*);
- Загрузка файла описания нейронной сети (*xml*);
- Предварительная обработка и загрузка изображений на вычислитей модуль;
- Запуск обработки на вычислительном модуле;
- Обработка и вывод результатов.

При запуске программы появляется главное окно программы со строками меню, состояния и клиентской областью.

В меню расположены органы управления программой. В клиентской области располагается изображение и результаты обработки. В строке состояния выводится следующая информация:

- Выбранный ускоритель (симулятор, MC121.01, MC127.05, NMStick, NMCard);
- Выбранная модель нейронной сети;
- Выбранное описание нейронной сети;
- Выбранное изображение;
- Скорость обработки (кадров/с).

Внешний вид программы в процессе работы приведён на рисунке 7.1.





Рисунок 7.1 - Внешний вид программы в процессе работы



Выбор модуля осуществляется в диалоговом окне *Open Board*, внешний вид которого приведён на рисунке <u>7.2</u>. Окно содержит список доступных в системе ускорителей, кнопки выбора, а также возможность идентификации устройства посредством светодиодной индикации.



Рисунок 7.2 - Внешний вид окна выбора модуля

Выбор описания нейронной сети осуществляется в диалоговом окне *Open Description*. Описание нейронной сети осуществляется в формате *XML* и содержит информацию о необходимом формате изображения, а также параметры интерпретации выходных параметров.

Параметр	Описание
model	Строка с названием файла модели нейронной сети в одном из следующих форматах:
	• *. <i>nm7</i> - описание модели для загрузки на модули <i>MC121.01</i> и <i>NMStick</i> .
	• <i>*.nm8</i> - описание модели для загрузки на модули <i>MC127.05, NMCard</i> или на симулятор.
	 *.<i>pb</i> или *.<i>onnx</i> - описание модели в формате ONNX Protobuf. Модель может быть загружена на любой модуль. Перед загрузкой на модуль модель предварительно конвертируется в *.<i>nm7</i> или в *.<i>nm8</i> в зависимости от типа модуля.
	 *.cfg - описание модели в формате DARKNET. Файл с весовыми коэффициентами *.weights должен размещаться в том же каталоге и иметь то же название, что и файл описания модели *.cfg. Модель может быть загружена на любой модуль. Перед загрузкой на модуль модель предварительно конвертируется в *.nm7 или в *.nm8 в зависимости от типа модуля.



Параметр	Описание		
	Можно указать как абсолютный путь к файлу, так и путь, относительно размещения файла описания XML.		
format	Формат пикселя, к которому будет преобразованы пиксели входного изображения перед обработкой. Может принимать значения:		
	• rgb		
	• rbg		
	• grb		
	• gbr		
	• brg		
	• bgr		
	Это значение соответствует порядку следования каналов входного тензора.		
divider	Масштабирующий коэффициент для каждой цветовой компоненты пикселя входного изображения. Используется при преобразовании изображения во входной тензор (см. раздел <u>"Подготовка изображений"</u>).		
adder	Коэффициент смещения для каждой цветовой компоненты пикселя входного изображения. Используется при преобразовании изображения во входной тензор (см. раздел "Подготовка изображений").		
neural_network	Выбранная нейронная сеть Может принимать значения:		
	• <i>classifier</i> - нейронная сеть - классификатор, например <i>ALEXNET</i> , <i>RESNET</i> , <i>SQUEEZENET</i> .		
	• <i>unet</i> - нейронная сеть семейства U-Net.		
	• <i>yolo2</i> - нейронная сеть семейства YOLO V2.		
	• <i>yolo3</i> - нейронная сеть семейства YOLO V3.		
	Результатом работы классификатора является вектор вероятностей принадлежности изображения определённым классам. В демонстрационной программе классы и вероятности обнаружения выводятся во всплывающем окне.		
	Сети YOLO в результате обработки выдают информацию о нескольких обнаруженных объектах и их расположении на исходном изображении. В демонстрационной программе		



Параметр	Описание		
	обнаруженые объекты обозначаются непосредственно на исходном.изображении в описывающих прямоугольниках.		
yolo_anchors	Параметры начальной инициализации детектируемых прямоугольников (только для сетей YOLO2 и YOLO3).		
yolo_confidence_threshold	Порог для отображения результатов детектирования (только для сетей <i>YOLO2</i> и <i>YOLO3</i>).		
yolo_iou_threshold	Порог пересекающихся прямоугольников детектирования (только для сетей <i>YOLO2</i> и <i>YOLO3</i>)		
labels	Список строк с названиями классов, которые будут отображаться в окне результата обработки.		

Диалоговое окно *Open Picture* позволяет открыть одно или несколько изображений для обработки.

При наличии всех данных становится доступной кнопка Run, запускающая обработку. Запуск можно инициировать с помощью комбинации "Ctrl+R". Для каждого переданного изображения осуществляется предварительная обработка в соответствии с заданным описанием нейронной сети для последующей обработки на ускорителе. В результате обработки заполняется выходная структура с N тензорами. В зависимости от нейронной сети, выходная структура интерпретируется различным образом. Для сетей классификации (таких как SqueezeNet) появляется дополнительное окно с классами и вероятностью. Внешний вид этого окна приведён на рисунке 7.3.



Classification Results 😐 💷 💌			
	Label	Probability	
1	Pembroke	29.1992	
2	Cardigan	26.5502	
3	Chihuahua	24.7543	
4	dingo	24.1734	
5	Labrador_retriever	23.6046	
6	Eskimo_dog	23.4812	
7	Staffordshire_bull	22.9756	
8	golden_retriever	22.6645	
9	basenji	22.3377	
10	Pomeranian	21.7855	

Рисунок 7.3 - Внешний вид окна результатов классификации

Для сетей, которые выполняют местоопределение объекта (таких как YOLO), результат показывается непосредственно на исходном изображении. После окончания работы ускорителя осуществляется наложение прямоугольников на обнаруженные объекты, приводится название класса и вероятность обнаружения. Пример наложения результата на изображение приведён на рисунке <u>7.4</u>.





Рисунок 7.4 - Отображение результатов детектирования

Программа позволяет обработать последовательность изображений. Для этого необходимо при выборе изображений выделить несколько файлов. После этого при запуске обработки (*Run*) будет обработано одно изображение. При повторном запуске - второе и т. д. Можно также запустить обработку в автоматическом режиме (*Run Auto*). Для останова автоматической обработки нажмите клавишу "*пробел*".

Для отображения дополнительной информации используется строка состояния. В строку выводится название выбранного устройства, файла описания нейронной сети, файла изображения, а также измеренное программой значение fps с учётом пересылки кадра и полученрия результата. Скорость обработки без учёта времени пересылок выводится в скобках.



8. Пример использования NMDL

В примере демонстрируется применение библиотеки *NMDL*, программных модулей для компиляции модели и подготовки изображений. Пример состоит из файла исходного кода *example.cpp* и скрипта сборки *CMakeLists.txt* в каталоге "*examples*" установочной директории.

Предполагается, что исходные данные для обработки находятся в каталоге "nmdl_ref_data/squeezenet" в каталоге установки, это означает, что в примере демонстрируется обработка нейронной сети squeezenet. Исходными данными являются:

- Модель нейронной сети в формате ONNX файл "nmdl ref data/squeezenet/model.pb"
- Обрабатываемое изображение в формате ВМР файл "nmdl_ref_data/squeezenet/ frame.bmp"

В примере показаны вызовы функций библиотек в порядке, необходимом для осуществления корректной работы.

```
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 = {
                                                            "OK"},
     {NMDL COMPILER RESULT OK,
014
015 {NMDL_COMPILER_RESULT_MEMORY_ALLOCATION_ERROR, "MEMORY ALLOCATION ERROR"},
     {NMDL_COMPILER_RESULT_MODEL_LOADING_ERROR, "MODEL_LOADING_ERROR"
{NMDL_COMPILER_RESULT_INVALID_PARAMETER, "INVALID_PARAMETER"},
{NMDL_COMPILER_RESULT_INVALID_MODEL, "INVALID_MODEL"},
016
                                                            "MODEL LOADING ERROR" },
017
     {NMDL COMPILER RESULT INVALID MODEL,
                                                           "INVALID MODEL" },
018
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"}
    {NMDL RESULT INVALID FUNC PARAMETER,
                                                        "INVALID FUNC PARAMETER" },
031
                                                        "NO LOAD LIBRARY" },
032 {NMDL RESULT NO LOAD LIBRARY,
     {NMDL_RESULT_NO_BOARD,
{NMDL_RESULT_BOARD_RESET_ERROR,
033
                                                        "NO BOARD" },

      034
      {NMDL_RESULT_BOARD_RESET_ERROR,
      "BOARD_RESET_ERROR"},

      035
      {NMDL_RESULT_INIT_CODE_LOADING_ERROR,
      "INIT_CODE_LOADING_ERROR"},

     {NMDL RESULT CORE HANDLE RETRIEVAL ERROR, "CORE HANDLE RETRIEVAL ERROR"},
036
      {NMDL RESULT FILE LOADING ERROR,
                                                        "FILE LOADING ERROR" },
037
                                                        "MEMORY WRITE ERROR" },
       {NMDL RESULT MEMORY WRITE ERROR,
038
                                                        "MEMORY_READ_ERROR" },
039
     {NMDL_RESULT_MEMORY_READ_ERROR,
040
      {NMDL RESULT MEMORY ALLOCATION ERROR,
                                                        "MEMORY ALLOCATION ERROR" },
                                                        "MODEL_LOADING_ERROR" },
041
       {NMDL RESULT MODEL LOADING ERROR,
042
      {NMDL RESULT INVALID MODEL,
                                                        "INVALID MODEL"},
      {NMDL_RESULT_BOARD_SYNC_ERROR,
                                                         "BOARD SYNC ERROR" },
043
044
       {NMDL_RESULT_BOARD_MEMORY_ALLOCATION_ERROR, "BOARD_MEMORY_ALLOCATION_ERROR"},
```



```
{NMDL RESULT NN CREATION ERROR,
045
                                                    "NN CREATION ERROR" },
046
      {NMDL RESULT NN LOADING ERROR,
                                                    "NN LOADING ERROR" },
047
      {NMDL RESULT NN INFO RETRIEVAL ERROR,
                                                    "NN INFO RETRIEVAL ERROR" },
                                                    "MODEL IS TOO BIG" },
      {NMDL RESULT MODEL IS TOO BIG,
048
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) {
     throw std::runtime error(function name + ": " + map[result]);
054
055 }
056
        return NMDL RESULT OK;
057 }
0.5.8
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 = Ou;
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;</pre>
104 for(std::size t i = 0; i < model info.output tensor num; ++i) {
    std::cout << "Output tensor " << i << ": " <<
105
       model_info.output_tensors[0].width << ", " <</pre>
106
107
      model info.output tensors[0].height << ", " <</pre>
108
     model_info.output_tensors[0].depth <<</pre>
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,
```



```
std::uint32 t color format, float divider, float adder) {
116
117 auto bmp frame = ReadFile<char>(filename);
118
    std::vector<float> nm frame(NMDL IMAGE CONVERTER RequiredSize(
      width, height, board_type));
119
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("Imnage conversion error");
123 }
124 return nm frame;
125 }
126
127 void WaitForOutput (NMDL HANDLE nmdl, std::uint32 t batch num,
     std::vector<float> &output) {
128
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) {
     std::cout << "FPS: " << fps << std::endl;</pre>
136
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
    BOARD TYPE == NMDL BOARD TYPE MC12101 ?
151
     NMDL_COMPILER_BOARD_TYPE_MC12101 :
152
    NMDL_COMPILER_BOARD_TYPE_MC12705;
153
154 const std::uint32 t IMAGE CONVERTER BOARD TYPE =
    BOARD_TYPE == NMDL_BOARD_TYPE_MC12101 ?
155
156
     NMDL IMAGE CONVERTER BOARD TYPE MC12101 :
     NMDL IMAGE CONVERTER BOARD TYPE MC12705;
157
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;</pre>
178
     ShowNMDLVersion();
179
180
      std::cout << "Board detection... " << std::endl;</pre>
181
      CheckBoard (BOARD TYPE);
182
      std::cout << "Compile model... " << std::endl;</pre>
183
184
      auto nm model = CompileModel(ONNX MODEL FILENAME, COMPILER BOARD TYPE);
185
      std::cout << "NMDL initialization... " << std::endl;</pre>
186
```



```
187
    Call(NMDL Create(&nmdl), "Create");
     Call(NMDL_Initialize(nmdl, BOARD_TYPE, 0, nm_model.data(),
188
189
       nm model.size(), 0), "Initialize");
190
191
    std::cout << "Get model information... " << std::endl;</pre>
192
     auto model info = GetModelInformation(nmdl);
193
194
     std::cout << "Prepare frame... " << std::endl;</pre>
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
     std::cout << "Process single frames... " << std::endl;</pre>
201
202
     for(std::size_t i = 0; i < FRAMES; ++i) {</pre>
       Call(NMDL ProcessFrame(nmdl, 0, frame.data()), "ProcessFrame");
203
204
      WaitForOutput(nmdl, 0, output);
205
     }
206
     NMDL Release(nmdl);
207
    std::cout << "Batch mode process frames... " << std::endl;</pre>
208
    std::uint32_t cnt_in = 0;
209
210
     std::uint32 t cnt out = 0;
211 Call(NMDL Initialize(nmdl, BOARD TYPE, 0, nm model.data(),
       nm_model.size(), 1), "Initialize");
212
213
     for(auto i = 0u; i < BATCHES; ++i) {</pre>
      Call(NMDL ProcessFrame(nmdl, (cnt in++) % BATCHES,
214
        frame.data()), "ProcessFrame");
215
216
217
     for(auto i = BATCHES; i < FRAMES; ++i) {</pre>
218
    WaitForOutput(nmdl, (cnt out++) % BATCHES, output);
219
     Call(NMDL ProcessFrame(nmdl, (cnt in++) % BATCHES,
220
        frame.data()), "ProcessFrame");
221
222 for(auto i = Ou; i < BATCHES; ++i) {
     WaitForOutput(nmdl, (cnt_out++) % BATCHES, output);
223
224
      }
225 }
226 catch (std::exception& e) {
227
     std::cerr << e.what() << std::endl;</pre>
228 }
229 NMDL Release(nmdl);
230 NMDL Destroy(nmdl);
231
232
       return 0;
233 }
```

Здесь выполняются следующие действия:

1.. 8: Подключение заголовочных файлов. "*nmdl.h*" - описание библиотеки нейросетевой обработки, "*nmdl_compiler.h*" - описание библиотеки для компиляции моделей, "*nmdl_image_converter.h*" - описание библиотеки для подготовки изображений.

12 .. 26: Функция-обёртка для вызовов функций библиотеки компиляции моделей. В случае ошибки формирует вывод об ошибке и инициирует исключение.

28.. 57: Функция-обёртка для вызовов функций библиотеки нейросетевой обработки. В случае ошибки формирует вывод об ошибке и инициирует исключение.

59... 70: Универсальная функция для чтения данных из файла в вектор.

72... 77: Функция вывода версии NMDL <u>NMDL_GetLibVersion</u>.



79... 87: Функция проверки наличия модуля заданного типа. Фактически производится запрос количества обнаруженных модулей заданного типа - <u>NMDL_GetBoardCount</u>.

89 .. 98: Функция компиляции исходной модели. Вызывается функция <u>NMDL_COMPILER_CompileONNX</u>. Здесь происходит выделение памяти внутри функции компиляции, поэтому после работы выделенная память освобождается вызовом <u>NMDL_COMPILER_FreeModel</u>, а результат компиляции копируется в возвращаемый вектор float. В целевой программе можно выполнить предварительную компиляцию модели с помощью утилиты *nmdl_compiler_console* так, как описано в разделе "Компиляция модели".

100.. *112*: Функция получения и вывода информации о параметрах выходных тензоров. Используется вызов <u>NMDL_GetModelInfo</u>.

114 .. 125: Функция подготовки кадра. Здесь выполняется чтение изображения из файла, его декодирование и предобработка (<u>NMDL_IMAGE_CONVERTER_Convert</u>). Описание предобработки приводится в разделе <u>"Подготовка изображений"</u>.

127 .. 142: Функция ожидания обработки кадра. Принимает номер кластера на котором производится обработка (*batch_num*) и буфер-вектор для хранения результата обработки. Функция блокируется в цикле запроса статуса обработки (<u>NMDL_GetStatus</u>). После получения статуса *NMDL_PROCESS_FRAME_STATUS_FREE* производится копирование результата вызовом <u>NMDL_GetOutput</u>.

147: Задание типа модуля ускорителя. В зависимости от типа ускорителя определяются и типы для работы с библиотекой компиляции моделей и библиотекой подготовки изображений. В примере используется симулятор модуля MC127.05. Для работы с другими типами раскомментируйте соответствующую строку.

147: Setting the type of the accelerator module. Depending on the type of accelerator, the types for working with the model compilation library and the image preparation library are also defined. The example uses the MC127.05 module simulator. For other types, uncomment the appropriate line.

158.. 164: Задаются имена файлов с исходными моделью и изображением.

166 ... 169: Задаются параметры для подготовки изображения. Для исходной модели требуется подготовить изображение с пикселями в формате blue-green-red. Каждый пиксель модифицируется по формуле Y = X / 1.0 - 114. Это преобразование необходимо для обработки модели *squeezenet*. Для других моделей необходимо использовать соответствующие параметры, которые определяются при создании модели и являются исходными данными, привязанными к конкретной модели. Подробнее о подготовке изображений см. в разделе <u>"Подготовка изображений"</u>.

171: Задаётся количество кластеров при пакетной обработке.

172: Задаётся количество обрабатываемых кадров. В примере производится многократная обработка одного кадра.

Далее в главной функции выполняется последовательный вызов описанных вспомогательных функций.



187 .. 189: Инициализация NMDL. Посредством вызова функции <u>NMDL_Create</u> осуществляется инициализация внутренних структур библиотеки NMDL. В функции <u>NMDL_Initialize</u> производится загрузка скомпилированной модели в выбранный ускоритель.

201 .. 205: Пример последовательной обработки кадров в режиме разделения данных по кластерам (см. раздел <u>"Режимы обработки"</u>).

208 .. 224: Пример последовательной обработки кадров в режиме пакетной обработки (см. раздел <u>"Режимы обработки"</u>).

229 .. 230: При завершении работы освобождаются выделенные ресурсы (<u>NMDL_Release</u> и <u>NMDL_Destroy</u>).



9. Описание идентификаторов, функций и структур NMDL

9.1. Идентификаторы и структуры

9.1.1. NMDL_BOARD_TYPE

Типы модулей.

```
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 симулятор модуля MC127.05.
- *NMDL_BOARD_TYPE_MC12101* модуль *MC121.01*.
- *NMDL_BOARD_TYPE_MC12705* модуль *MC127.05*.
- NMDL_BOARD_TYPE_NMSTICK модуль NMStick.
- *NMDL_BOARD_TYPE_NMCARD* модуль *NMCard*.

9.1.2. NMDL_ModelInfo

Структура с информацией о модели.

```
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* описание входного тензора (входное изображение),
- *output_tensor_num* количество выходных тензоров,
- *output_tensors* выходные тензоры.

NMDL_MAX_OUTPUT_TENSORS - максимальное количество выходных тензоров.

9.1.3. NMDL_PROCESS_FRAME_STATUS

Идентификатор статуса обработки кадра.

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 обработка кадра не производится,
- *NMDL_PROCESS_FRAME_STATUS_BUSY* выполняется обработка кадра.

9.1.4. NMDL_RESULT

Возвращаемый результат.

```
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* нет ошибок,
- NMDL_RESULT_INVALID_FUNC_PARAMETER неверный параметр,
- *NMDL_RESULT_NO_BOARD* нет модуля,
- *NMDL_RESULT_BOARD_RESET_ERROR* ошибка сброса модуля,
- *NMDL_RESULT_INIT_CODE_LOADING_ERROR* ошибка загрузки кода инициализации библиотеки загрузки и обмена,
- *NMDL_RESULT_CORE_HANDLE_RETRIEVAL_ERROR* ошибка получения идентификатора вычислителя,
- *NMDL_RESULT_FILE_LOADING_ERROR* ошибка загрузки программного образа,
- *NMDL_RESULT_MEMORY_WRITE_ERROR* ошибка записи в память модуля,
- NMDL_RESULT_MEMORY_READ_ERROR ошибка чтения из памяти модуля,
- NMDL RESULT MEMORY ALLOCATION ERROR ошибка выделения памяти,



- NMDL_RESULT_MODEL_LOADING_ERROR ошибка загрузки модели нейронной сети,
- *NMDL_RESULT_INVALID_MODEL* ошибка в модели,
- *NMDL_RESULT_BOARD_SYNC_ERROR* ошибка синхронизации с модулем,
- *NMDL_RESULT_BOARD_MEMORY_ALLOCATION_ERROR* ошибка выделения памяти на модуле,
- NMDL_RESULT_NN_CREATION_ERROR ошибка создания модели на модуле,
- *NMDL_RESULT_NN_LOADING_ERROR* ошибка загрузки модели нейронной сети,
- *NMDL_RESULT_NN_INFO_RETRIEVAL_ERROR* ошибка запроса информации о модели,
- *NMDL_RESULT_MODEL_IS_TOO_BIG* модель не может быть размещена в памяти модуля,
- *NMDL_RESULT_NOT_INITIALIZED* библиотека функций NMDL не инициализирована,
- *NMDL_RESULT_BUSY* устройство находится в состоянии обработки кадра,
- *NMDL_RESULT_UNKNOWN_ERROR* неизвестная ошибка.

9.1.5. NMDL_Tensor

Структура, описывающая тензор.

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

- *width* ширина тензора,
- *height* высота тензора,
- *depth* глубина тензора.

9.2. Функции

9.2.1. NMDL_Blink

Светодиодная индикация для идентификации модуля.

```
NMDL_RESULT NMDL_Blink(
    unsigned int board_type,
```



```
unsigned int board_number
);
```

)

- *board_type* [in] тип модуля, на котором вызвается процедура светодиодной индикации. Одно из значений перечисления <u>NMDL_BOARD_TYPE</u>.
- *board_number* [in] порядковый номер модуля, на котором вызвается процедура светодиодной индикации.

9.2.2. NMDL_Create

Создание экземпляра NMDL и получение идентификатора экземпляра NMDL.

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

• *nmdl* - [out] идентификатор экземпляра *NMDL*.

После работы с экземпляром NMDL необходимо освободить выделенные ресурсы вызовом <u>NMDL_Destroy</u>.

9.2.3. NMDL_Destroy

Удаление экземпляра NMDL.

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

• *nmdl* - [in] идентификатор экземпляра *NMDL*.

Функция вызывается для освобождения ресурсов, выделенных при вызовах <u>NMDL_Create</u> и<u>NMDL_Initialize</u>.

9.2.4. NMDL_GetBoardCount

Запрос количества обнаруженных модулей заданного типа.

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

- *board_type* [in] тип опрашиваемых модулей. Одно из значений перечисления <u>NMDL BOARD TYPE</u>.
- *boards* [out] количество обнаруженных модулей.

Для симулятора MC127.05 (тип *NMDL_BOARD_TYPE_SIMULATOR*) всегда обнаруживается один модуль.

9.2.5. NMDL_GetLibVersion

Запрос версии библиотеки NMDL.

```
NMDL_RESULT NMDL_GetLibVersion(
    unsigned int *major,
    unsigned int *minor,
    unsigned int *patch
);
```

- *major* [out] старший номер версии.
- *minor* [out] младший номер версии.
- *patch* [out] номер патча.

9.2.6. NMDL_GetModelInfo

Запрос информации о модели.

```
NMDL_RESULT NMDL_GetModelInfo(
    NMDL_HANDLE nmdl,
    <u>NMDL_ModelInfo</u> *model_info
);
```

- *nmdl* [in] дескриптор *NMDL*
- *model_info* [out] информация о модели.

9.2.7. NMDL_GetOutput

Запрос результата обработки.

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

- *nmdl* [in] идентификатор экземпляра *NMDL*.
- *batch_num* [in] номер кластера для которого запрашивается результат обработки. Имеет значение только при работе с модулями *MC127.05*, *NMCard* или симулятором. При работе с *MC121.01* и *NMStick* необходимо установить 0.
- *output* [out] указатель на буфер выходных тензоров. Выходные тензоры размещаются друг за другом в порядке их появления в графе обработки, то есть упорядоченные по уровням в графе. Параметры тензоров определяются в структуре <u>NMDL_ModelInfo</u>. Размер требуемого буфера возвращается в функции <u>NMDL_RequiredOutputFloats</u>. Может принимать значение 0.



• *fps* - [out] производительность обработки (кадров в секунду). Значение производительности определено только для кластера 0 (batch_num = 0), для других кластеров записывает 0.0f. Может принимать значение 0.

9.2.8. NMDL_GetStatus

Запрос статуса обработки. Функция вызывается для проверки окончания обработки кадра.

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

- nmdl [in] идентификатор экземпляра NMDL
- batch_num [in] номер кластера для которого запрашивается статус обработки. Имеет значение только при работе с модулями MC127.05, NMCard или симулятором. При работе с MC121.01 и NMStick необходимо установить 0.
- *status* [out] состояное кластера в момент вызова функции. Одно из значений перечисления <u>NMDL_PROCESS_FRAME_STATUS</u>.

9.2.9. NMDL_Initialize

Инициализация NMDL.

```
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] идентификатор экземпляра *NMDL*.
- *board_type* [in] тип инициализируемого модуля. Одно из значений перечисления <u>NMDL_BOARD_TYPE</u>.
- board_number [in] порядковый номер инициализируемого модуля.
- *model* [in] указатель на буфер, содержащий модель в формате *nm7* (для *MC121.01* и *NMStick*) или *nm8* (для *MC127.05*, *NMCard* и симулятора).
- model_floats [in] размер модели в вещественных числах с одинарной точностью.
- *use_batch_mode* [in] флаг использования режима пакетной обработки. Имеет значение только для модулей *MC127.05*, *NMCard* или симулятора. При работе с *MC121.01* и *NMStick* необходимо установить 0.



После работы с функциями NMDL необходимо освободить выделенные ресурсы вызовом функции деинициализации <u>NMDL_Release</u>.

9.2.10. NMDL_ProcessFrame

Обработка одиночного кадра.

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

- *nmdl* [in] идентификатор экземпляра *NMDL*
- *batch_num* [in] номер кластера на котором запускается обработка. Имеет значение только при работе с модулями *MC127.05*, *NMCard* или симулятором. При работе с *MC121.01* и *NMStick* необходимо установить 0.
- *frame* [in] указатель на буфер обрабатываемого кадра.

9.2.11. NMDL_Release

Деинициализация NMDL.

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

• *nmdl* - [in] идентификатор экземпляра *NMDL*

9.2.12. NMDL_RequiredOutputFloats

Возвращает размер буфера в элементах *float32*, необходимый для размещения результата (*output*) при вызове <u>NMDL_GetOutput</u>.

```
unsigned int NMDL_RequiredOutputFloats(
    const <u>NMDL_ModelInfo</u> *model_info
);
```

• *model_info* - [out] информация о модели, ренее полученная с помощью вызова <u>NMDL_GetModelInfo</u>.

10. Описание идентификаторов и функций nmdl_compiler

10.1. Идентификаторы

10.1.1. NMDL_COMPILER_BOARD_TYPE

Типы модулей.

```
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 модули MC121.01 и NMStick.
- *NMDL_COMPILER_BOARD_TYPE_MC12705* модули *MC127.05*, *NMCard* и симулятор.

10.1.2. NMDL_COMPILER_RESULT

Возвращаемый результат.

```
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* нет ошибок,
- *NMDL_COMPILER_RESULT_MEMORY_ALLOCATION_ERROR* ошибка выделения памяти,
- *NMDL_COMPILER_RESULT_MODEL_LOADING_ERROR* ошибка загрузки модели нейронной сети,
- NMDL_COMPILER_RESULT_INVALID_PARAMETER неверный параметр,
- *NMDL_COMPILER_RESULT_INVALID_MODEL* ошибка в модели,
- *NMDL_COMPILER_RESULT_UNSUPPORTED_OPERATION* модель содержит неподдерживаемую операцию.



10.2. Функции

10.2.1. NMDL_COMPILER_CompileDarkNet

Компиляция исходной модели в формате DarkNet.

```
NMDL_COMPILER_RESULT NMDL_COMPILER_CompileDarkNet(
    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] идентификатор типа модуля для которого производится компиляция модели. Одно из значений перечисления <u>NMDL_COMPILER_BOARD_TYPE</u>.
- *src_model* [in] буфер исходной модели в формате *DarkNet*, предварительно считанной из файла .*cfg*.
- *src_model_size* [in] размер буфера исходной модели в байтах.
- src_weights [in] буфер коэффициентов, предварительно считанный из файла . weights.
- *src_weights_size* [in] размер буфера коэффициентов в байтах.
- *dst_model* [out] буфер компилированной модели. Выделение памяти происходит в функции.
- *dst_model_floats* [out] размер буфера компилированной модели в *float32*.

10.2.2. NMDL_COMPILER_CompileONNX

Компиляция исходной модели в формате ONNX.

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

- *board* [in] идентификатор типа модуля для которого производится компиляция модели. Одно из значений перечисления <u>NMDL_COMPILER_BOARD_TYPE</u>.
- *src_model* [in] буфер исходной модели в формате *ONNX*, предварительно считанной из файла *.pb*.
- src_model_size [in] размер буфера исходной модели в байтах.





- *dst_model* [out] буфер компилированной модели. Выделение памяти происходит в функции.
- *dst_model_floats* [out] размер буфера компилированной модели в *float32*.

10.2.3. NMDL_COMPILER_FreeModel

Освобождение выделенной при компиляции памяти.

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

-);
- *board* [in] идентификатор типа модуля для которого производедена компиляция модели. Одно из значений перечисления <u>NMDL_COMPILER_BOARD_TYPE</u>.
- dst model -[in] буфер освобождаемой Память буфера памяти. для NMDL COMPILER CompileDarkNet быля выделена при вызове или NMDL COMPILER CompileONNX.

10.2.4. NMDL_COMPILER_GetLastError

Возвращает константную строку с описанием последней ошибки.

const char *NMDL_COMPILER_GetLastError();



11. Описание идентификаторов и функций nmdl_image_converter

11.1. Идентификаторы и структуры 11.1.1. NMDL_IMAGE_CONVERTER_BOARD_TYPE

Типы модулей.

```
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* модули *MC121.01* и *NMStick*.
- *NMDL_IMAGE_CONVERTER_BOARD_TYPE_MC12705* модули *MC127.05*, *NMCard* и симулятор.

11.1.2. NMDL_IMAGE_CONVERTER_COLOR_FORMAT

Идентификаторы формата пикселя. Описывает порядок следования цветовых компонент RGB в пикселе.

```
typedef enum tagNMDL_IMAGE_CONVERTER_COLOR_FORMAT {
    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_BT,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_INTENSITY,
    NMDL_IMAGE_CONVERTER_COLOR_FORMAT_INTENSITY,
} NMDL_IMAGE_CONVERTER_COLOR_FORMAT;
```

11.2. Функции

11.2.1. NMDL_IMAGE_CONVERTER_Convert

Подготовка изображения.

```
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] буфер с образом исходного изображения. Содержимое буфера соответствует содержимому файла изображения. Изображения могут быть представлены в форматах .bmp, .gif, .jpg и .png.
- *dst* [out] буфер для подготовленного изображения. Память для буфера должна быть предварительно выделена. Размер буфера должен быть не менше значения, возвращаемого функцией <u>NMDL_IMAGE_CONVERTER_RequiredSize</u>.
- *src_size* [in] размер буфера исходного изображения в байтах.
- *dst_width* [in] ширина подготовленногоизображения.
- *dst_height* [in] высота подготовленногоизображения.
- *dst_color_format* [in] идентификатор формата пикселя подготовленного изображения. Одно из значений перечисления <u>NMDL_IMAGE_CONVERTER_COLOR_FORMAT</u>.
- *rgb_divider* [in] делитель в выражении dst = src / divider+ adder. Приведённая операция выполняется над каналами каждого пикселя исходного изображения.
 rgb_divider[0] для красного канала, rgb_divider[1] для зелёного канала, rgb_divider[2] для голубого канала.
- *rgb_adder* [in] слагаемое в выражении dst = src / divider+ adder. Приведённая операция выполняется над каналами каждого пикселя исходного изображения. rgb_adder[0] для красного канала, rgb_adder[1] для зелёного канала, rgb_adder[2] для голубого канала.
- *board_type* [in] идентификатор типа вычислительного модуля на котором предполагается обработка. Одно из значений перечисления <u>NMDL_IMAGE_CONVERTER_BOARD_TYPE</u>.

Для изображений с градацией серого, когда поле dst_color_format = NMDL_IMAGE_CONVERTER_COLOR_FORMAT_INTENSITY, используются только делитель rgb_divider[0] и слагаемое rgb_adder[0]. Остальные делители (rgb_divider[1] и rgb_divider[2]) и слагаемые (rgb_adder[1] и rgb_adder[2]) игнорируются и могут быть установлены в любые значения.

Возвращаемое значение: 0 - нормальное завершение, -1 - ошибка.

11.2.2. NMDL_IMAGE_CONVERTER_RequiredSize

Возвращает размер буфера в элементах *float32* для хранения подготовленного изображения.

```
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] ширина подготовленногоизображения.
- dst height [in] высота подготовленногоизображения.
- *dst_color_format* [in] идентификатор формата пикселя подготовленного изображения. Одно из значений перечисления <u>NMDL_IMAGE_CONVERTER_COLOR_FORMAT</u>.
- board type [in] идентификатор ٠ типа вычислительного модуля на предполагается обработка. котором Одно ИЗ значений перечисления NMDL IMAGE CONVERTER BOARD TYPE.

Возвращаемое значение: 0 - нормальное завершение, -1 - ошибка.





АО НТЦ "Модуль" А/Я 166, Москва, 125190, Россия Тел: +7 (499) 152-9698 Факс: +7 (499) 152-4661 E-Mail: rusales@module.ru WWW: http://www.module.ru

Напечатано в России. Дата публикации: 20/12/2021

©АО НТЦ "Модуль", 2021 Все права защищены.

Весе права часть информации, приведенная в данном документе, не может быть адаптирована или воспроизведена, кроме как согласно письменному разрешению владельцев авторских прав. АО НТЦ "Модуль" оставляет за собой право производить изменения как в описании, так и в самом продукте без дополнительных уведомлений. АО НТЦ "Модуль" не несет ответственности за любой ущерб, причиненный использованием информации в данном описании, ошибками или недосказанностью в описании, а также путем неправильного исползования продукта.

