



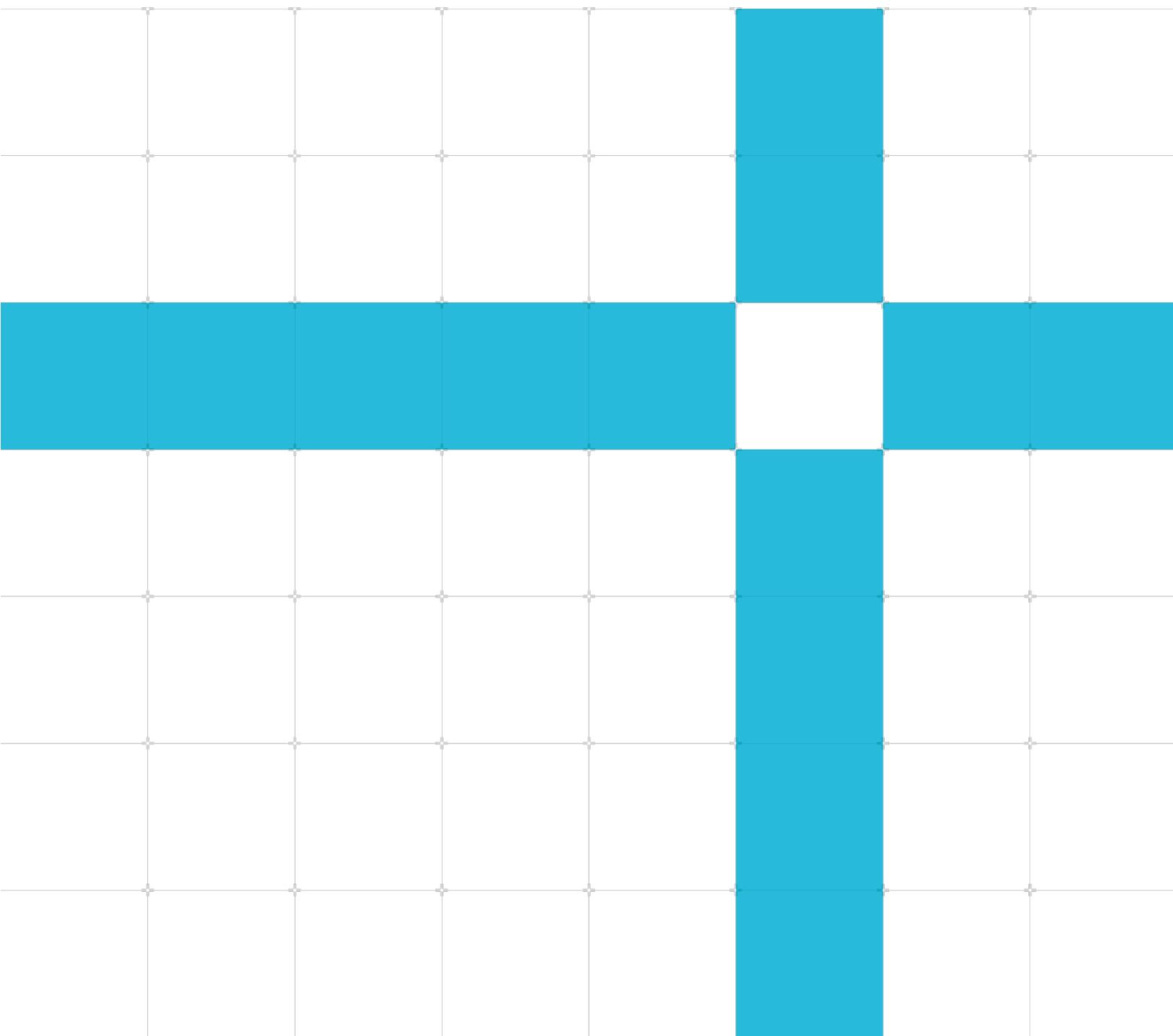
Object recognition with Arm NN and Raspberry Pi

Non-Confidential

Copyright © 2020 Arm Limited (or its affiliates).
All rights reserved.

Issue 0100

102274



Object recognition with Arm NN and Raspberry Pi

Copyright © 2020 Arm Limited (or its affiliates). All rights reserved.

Release information

Document history

Issue	Date	Confidentiality	Change
01	14 October 2020	Non-Confidential	First release

Non-Confidential Proprietary Notice

This document is protected by copyright and other related rights and the practice or implementation of the information contained in this document may be protected by one or more patents or pending patent applications. No part of this document may be reproduced in any form by any means without the express prior written permission of Arm. No license, express or implied, by estoppel or otherwise to any intellectual property rights is granted by this document unless specifically stated.

Your access to the information in this document is conditional upon your acceptance that you will not use or permit others to use the information for the purposes of determining whether implementations infringe any third party patents.

THIS DOCUMENT IS PROVIDED "AS IS". ARM PROVIDES NO REPRESENTATIONS AND NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, INCLUDING, WITHOUT LIMITATION, THE IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, NON-INFRINGEMENT OR FITNESS FOR A PARTICULAR PURPOSE WITH RESPECT TO THE DOCUMENT. For the avoidance of doubt, Arm makes no representation with respect to, and has undertaken no analysis to identify or understand the scope and content of, patents, copyrights, trade secrets, or other rights.

This document may include technical inaccuracies or typographical errors.

TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL ARM BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF ARM HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

This document consists solely of commercial items. You shall be responsible for ensuring that any use, duplication or disclosure of this document complies fully with any relevant export laws and regulations to assure that this document or any portion thereof is not exported, directly or indirectly, in violation of such export laws. Use of the word "partner" in reference to Arm's customers is not intended to create or refer to any partnership relationship with any other company. Arm may make changes to this document at any time and without notice.

If any of the provisions contained in these terms conflict with any of the provisions of any click through or signed written agreement covering this document with Arm, then the click through or signed written agreement prevails over and supersedes the conflicting provisions of these terms. This document may be translated into other languages for convenience, and you agree that if there is any conflict between the English version of this document and any translation, the terms of the English version of the Agreement shall prevail.

The Arm corporate logo and words marked with ® or ™ are registered trademarks or trademarks of Arm Limited (or its affiliates) in the US and/or elsewhere. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners. Please follow Arm's trademark usage guidelines at <http://www.arm.com/company/policies/trademarks>.

Copyright © 2020 Arm Limited (or its affiliates). All rights reserved.

Arm Limited. Company 02557590 registered in England.

110 Fulbourn Road, Cambridge, England CB1 9NJ.

(LES-PRE-20349)

Confidentiality Status

This document is Non-Confidential. The right to use, copy and disclose this document may be subject to license restrictions in accordance with the terms of the agreement entered into by Arm and the party that Arm delivered this document to.

Unrestricted Access is an Arm internal classification.

Web Address

www.arm.com

Contents

1 Overview	6
2 Before you begin	7
3 What is object detection?	10
4 Object detection model structure	12
4.1 YOLO and MobileNet-SSD.....	12
4.2 Region-Based Convolutional Neural Network.....	13
5 PyArmNN object recognition	14
5.1 Read from the video source.....	15
5.2 Prepare labels and model-specific functions.....	15
5.3 Create the parser and import a graph.....	15
5.4 Optimize the graph for the compute device.....	16
5.5 Create input and output binding information.....	16
5.6 Preprocess the captured frame.....	17
5.7 Make the input and output tensors.....	17
5.8 Execute inference.....	17
5.9 Decode and process the inference output.....	17
5.10 Draw the bounding boxes.....	18
5.11 Run the application.....	18
6 Arm NN C++ API object recognition	19
6.1 Read from the video source.....	19
6.2 Prepare labels and model-specific functions.....	20
6.3 Create a network.....	20
6.4 Create the parser and import a graph.....	20
6.5 Optimize the graph for the compute device.....	20
6.6 Create input and output binding information.....	21
6.7 Object detection pipeline.....	22
6.8 Preprocess the captured frame.....	22
6.9 Execute inference.....	22
6.10 Decode and process the inference output.....	23

6.11 Draw the bounding boxes23

6.12 Run the application23

7 Related information 25

8 Next steps 26

1 Overview

This guide shows you how to use Arm NN and PyArmNN to build and run a real-time object detection system.

The system runs on a Raspberry Pi 4 with Raspbian 10 operating system.

This guide examines the following sample applications that ship as part of Arm NN and PyArmNN:

- [PyArmNN Object Detection Sample Application](#)
- [Arm NN Object Detection Example](#)

These sample applications take a model and video file or camera feed as input. The applications then run inference on each frame. Finally, the applications draw bounding boxes around detected objects, with the corresponding labels and confidence scores overlaid.

By understanding how these sample applications perform object detection, you can learn to write your own Machine Learning applications using Arm NN and PyArmNN.

2 Before you begin

To work through this guide, you need a Raspberry Pi 4 with Raspbian 10 operating system.

You must install the following software on the Raspberry Pi:

- Arm NN version 20.05 or newer.

See [armnn/README.md](#) for installation instructions.

- PyArmNN

See [pyarmnn/README.md](#) for installation instructions.

These software packages might have further prerequisites. See the installation instructions for each package for more information.

For complete and up-to-date installation information, always refer to the previous README links. However, for your convenience, we provide, at the end of this section, a list of the commands that we used to install Arm NN and PyArmNN.

[PyArmNN object recognition](#) uses the `run_video_file.py` sample application. Follow the instructions in [PyArmNN Object Detection Sample Application](#) to install and build `run_video_file.py`.

[Arm NN C++ API object recognition](#) uses the Arm NN `object_detection_example` executable. Follow the instructions in [Object Detection Example](#) to install and build `object_detection_example`.

The installation links referenced earlier use OpenCV 3. If you want to use OpenCV 4 instead, you can follow the instructions in [Install OpenCV 4 on your Raspberry Pi](#). The sample application in this guide needs to have the `WITH_GTK` and `WITH_OPENCL` flags enabled. In [Step #5: CMake and compile OpenCV 4 for your Raspberry Pi](#), enable the flags using the following command:

```
cmake -DCMAKE_BUILD_TYPE=RELEASE -DCMAKE_INSTALL_PREFIX=/usr/local -DENABLE_NEON=ON
-DOPENCV_ENABLE_NONFREE=ON -DINSTALL_PYTHON_EXAMPLES=OFF -DBUILD_EXAMPLES=OFF
-DBUILD_DOCS=OFF -DBUILD_TEST=OFF -DWITH_GTK=ON -DWITH_OPENCL=ON -DWITH_LIBV4L=ON
-DWITH_V4L=OFF
```

Update your `PYTHONPATH` environment variable to include `/usr/local/python/cv2/python-<version>`.

The following code builds Arm NN and installs PyArmNN:

```
# Increase virtual memory swapfile allocation
sudo vi /etc/dphys-swapfile
# Find the following line:
#   CONF_SWAPSIZE=100
# Change this line to:
#   CONF_SWAPSIZE=1024
sudo /etc/init.d/dphys-swapfile stop
sudo /etc/init.d/dphys-swapfile start

# Install SCONS and CMAKE
sudo apt-get update
sudo apt-get install scons
sudo apt-get install cmake
```

```

# Install Arm NN prerequisites
mkdir armnn-tflite && cd armnn-tflite
export BASEDIR=`pwd`
git clone https://github.com/Arm-software/ComputeLibrary.git
git clone https://github.com/Arm-software/armnn
wget https://dl.bintray.com/boostorg/release/1.64.0/source/boost_1_64_0.tar.bz2
tar xf boost_1_64_0.tar.bz2
git clone -b v3.5.0 https://github.com/google/protobuf.git
git clone https://github.com/tensorflow/tensorflow.git
cd tensorflow/
git checkout 590d6eef7e91a6a7392c8ffffb7b58f2e0c8bc6b
git clone https://github.com/google/flatbuffers.git
cd $BASEDIR/ComputeLibrary
scons extra_cxx_flags="-fPIC" benchmark_tests=0 validation_tests=0 neon=1
cd $BASEDIR/boost_1_64_0
./bootstrap.sh
./b2 --build-dir=$BASEDIR/boost_1_64_0/build toolset=gcc link=static cxxflags=-fPIC --
with-filesystem --with-test --with-log --with-program_options install --
prefix=$BASEDIR/boost
cd $BASEDIR/protobuf
git submodule update --init --recursive
sudo apt-get install autoconf
sudo apt-get install libtool
./autogen.sh
./configure --prefix=$BASEDIR/protobuf-host
make
make install
cd $BASEDIR/tensorflow
../armnn/scripts/generate_tensorflow_protobuf.sh ../tensorflow-protobuf ../protobuf-
host
cd $BASEDIR
git clone https://github.com/google/flatbuffers.git
cd $BASEDIR/flatbuffers
cmake -G "Unix Makefiles" -DCMAKE_BUILD_TYPE=Release
make

#Install SWIG
sudo apt-get install libpcre3 libpcre3-dev
cd $BASEDIR
mkdir swig
cd swig
wget http://prdownloads.sourceforge.net/swig/swig-4.0.2.tar.gz
chmod 777 swig-4.0.2.tar.gz
tar -xzvf swig-4.0.2.tar.gz
cd swig-4.0.2/
./configure --prefix=/home/pi/armnn-tflite/swigtool/
sudo make
sudo make install
sudo vi /etc/profile

# Add the following lines to /etc/profile
# export SWIG_PATH=/home/pi/armnn-tflite/swigtool/bin
# export PATH=$SWIG_PATH:$PATH

source /etc/profile

# Build Arm NN
cd $BASEDIR/armnn
mkdir build
cd build
cmake .. -DARMCOMPUTE_ROOT=$BASEDIR/ComputeLibrary -
DARMCOMPUTE_BUILD_DIR=$BASEDIR/ComputeLibrary/build -DBOOST_ROOT=$BASEDIR/boost -
DTF_GENERATED_SOURCES=$BASEDIR/tensorflow-protobuf -DPROTOBUF_ROOT=$BASEDIR/protobuf-
host -DBUILD_TF_LITE_PARSER=1 -

```

```
DTF_LITE_GENERATED_PATH=$BASEDIR/tensorflow/tensorflow/lite/schema -
DFLATBUFFERS_ROOT=$BASEDIR/flatbuffers -
DFLATBUFFERS_LIBRARY=$BASEDIR/flatbuffers/libflatbuffers.a -DSAMPLE_DYNAMIC_BACKEND=1 -
DDYNAMIC_BACKEND_PATHS=$BASEDIR/armnn/src/dynamic/sample -DARMCOMPUTEONEON=1 -
DBUILD_TF_PARSER=1
make
cp $BASEDIR/armnn/build/*.so $BASEDIR/armnn/
cd /home/pi/armnn-tflite/armnn/src/dynamic/sample
mkdir build
cd build
cmake -DBOOST_ROOT=$BASEDIR/boost -
DBoost_SYSTEM_LIBRARY=$BASEDIR/boost/lib/libboost_system.a -
DBoost_FILESYSTEM_LIBRARY=$BASEDIR/boost/lib/libboost_filesystem.a -
DARMNN_PATH=$BASEDIR/armnn/libarmnn.so ..
make

# Install PYARMNN
# Following instructions for "Standalone build" from:
# https://git.mlplatform.org/ml/armnn.git/tree/python/pyarmnn/README.md

export SWIG_EXECUTABLE=$BASEDIR/swigtool/bin/swig
export ARMNN_INCLUDE=$BASEDIR/armnn/include/
export ARMNN_LIB=$BASEDIR/armnn/build/
cd $BASEDIR/armnn/python/pyarmnn
sudo apt-get install python3.6-dev build-essential checkinstall libreadline-gplv2-dev
libncursesw5-dev libssl-dev libsqlite3-dev tk-dev libgdbm-dev libc6-dev libbz2-dev

python3 setup.py clean --all
python3 swig_generate.py -v
python3 setup.py build_ext --inplace
python3 setup.py sdist
python3 setup.py bdist_wheel
pip3 install dist/pyarmnn-21.0.0-cp37-cp37m-linux_armv7l.whl
sudo pip3 install opencv-python==3.4.6.27
sudo apt-get install libblas-dev
sudo apt-get install libhdf5-dev
sudo apt-get install libhdf5-serial-dev
sudo apt-get install libatlas-base-dev
sudo apt-get install libjasper-dev
sudo apt-get install libqtgui4
sudo apt-get install libqt4-test
```

3 What is object detection?

Object detection is the process of identifying real-world objects, like people, cars, bottles, or sandwiches, in still images or videos. To understand an entire scene, a viewer must detect, identify, and locate multiple objects simultaneously.

Object detection is commonly used in applications like image retrieval, security, surveillance, and Advanced Driver Assistance Systems (ADAS).

There are several different approaches to performing object detection, including the following:

- Feature-based object detection
- Support Vector Machine (SVM) classification with Histogram of Oriented Gradient (HOG) feature descriptors
- Deep learning object detection

This guide focuses on deep learning object detection with Arm NN SDK.

Real-world use cases for object detection use cases include the following:

- Video surveillance

Video surveillance is a natural application of object detection techniques. The ability to identify and track multiple instances of an object in a scene is a key requirement of video surveillance systems. Granular object tracking provides valuable insights that can assist, for example, crowd security, analysis of store traffic, or stock movement on the factory floors.

The following image shows an example of object detection used in video surveillance:



Example of object detection for video surveillance

- Autonomous driving
Autonomous vehicles use real-time obstacle detection models to identify cars, pedestrians, bicycles, and traffic signs. The vehicle uses this information to decide whether to accelerate, decelerate, or turn.

The following image shows an example of real-time object detection used in an autonomous vehicle:

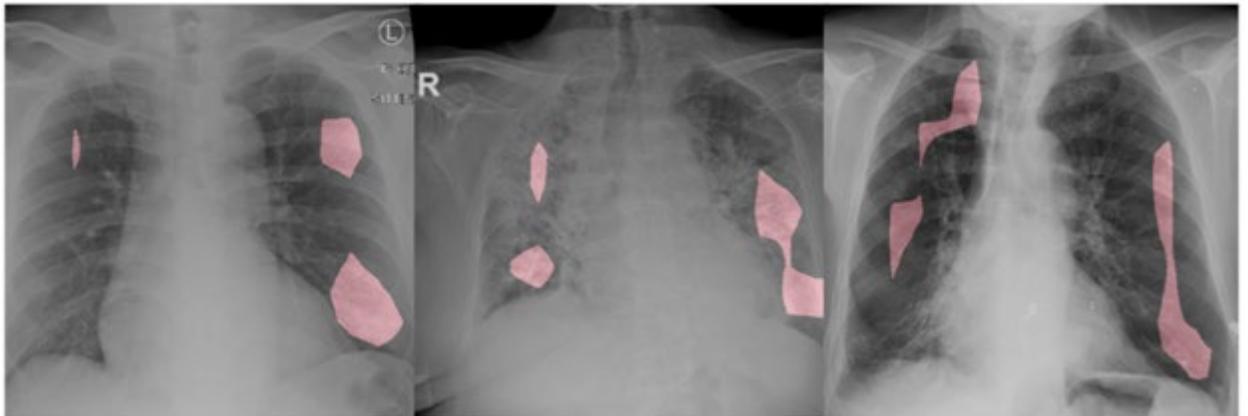


Example of object detection for autonomous vehicles (source: Mighty AI)

- Medical imaging

In health care, object detection assists doctors in diagnosis and treatment planning. For example, Arm ecosystem partner Darwin.ai developed COVID-Net, a deep neural network which examines chest X-rays to help medical professionals rapidly screen for COVID-19 infections.

The following image shows an example of object detection in medical X-rays:



Example of object detection in medical X-rays (source: DarwinAI)

- Manufacturing

Object detection helps to fully automate manufacturing systems. Object detection has applications in many manufacturing processes, including quality assurance, inventory management, and assembly line sorting.

4 Object detection model structure

Most deep learning-based object detection models have two parts:

- An encoder, which takes an image as input and runs it through a series of blocks and layers that extract features. The encoder then uses these features to locate and label objects.
- A decoder. Outputs from the encoder are then passed to a decoder. The decoder predicts bounding boxes and labels for each object.

The simplest decoder is a pure regressor. The regressor connects to the output of the encoder and directly predicts the location and size of each bounding box. The output of the model is the X, Y coordinate pair for the object and its extent. The disadvantage of using a pure regressor is that you must define the number of predicted objects ahead of time.

An extension of the regressor approach is a region proposal network. In this type of decoder, the model proposes regions of an image where it believes an object might reside. The pixels in these regions are fed into a classification network to determine a matching label. The region proposal network is a more accurate and flexible model that can process an arbitrary number of regions.

Single Shot Detectors (SSDs) seek to provide a middle ground between pure regressors and region proposal networks. Rather than using a subnetwork to propose regions, SSDs rely on a set of predetermined regions. A grid of anchor points is laid over the input image. At each anchor point, boxes of multiple shapes and sizes serve as regions.

For each box at each anchor point, the SSD model outputs:

- A prediction of whether an object exists within the region
- Modifications to the location and size of the box, to make the box fit the object more closely.

Because there are multiple boxes at each anchor point and anchor points might be close together, SSDs produce many potential detections that overlap. Post-processing must be applied to SSD outputs, to prune away most of these predictions, and pick the best one.

Object detectors output the location and label for each object. To benchmark the model performance, the most commonly-used metric is intersection-over-union (IOU). Given two bounding boxes, you compute the area of the intersection and divide by the area of the union. Metric values range from 0 (no interaction) to 1 (perfectly overlapping). For labels, you can use a simple percentage correct.

4.1 YOLO and MobileNet-SSD

Several models belong to the SSD family. The main differences between these variants are their encoders and the specific configuration of predetermined anchors.

YOLO v3 is a fast real-time object detection system. YOLO stands for You Only Look Once. Techniques like multi-scale predictions and improved backbone classifiers enable this fast performance. YOLO trains a single neural network model that directly predicts bounding boxes and class labels for each bounding box. You can find more details in the paper [YOLOv3: An Incremental Improvement](#).

MobileNet-SSD models feature a MobileNet-based encoder. SSDs are a good choice for models that are destined for mobile or embedded devices. For more information, see [MobileNetV2 + SSDLite with Core ML](#).

4.2 Region-Based Convolutional Neural Network

The Region-Based Convolutional Neural Network (R-CNN) family of methods do the following:

1. Generate candidate bounding boxes
2. Extract features from each candidate region using a deep convolutional neural network
3. Classify the features as one of the known classes

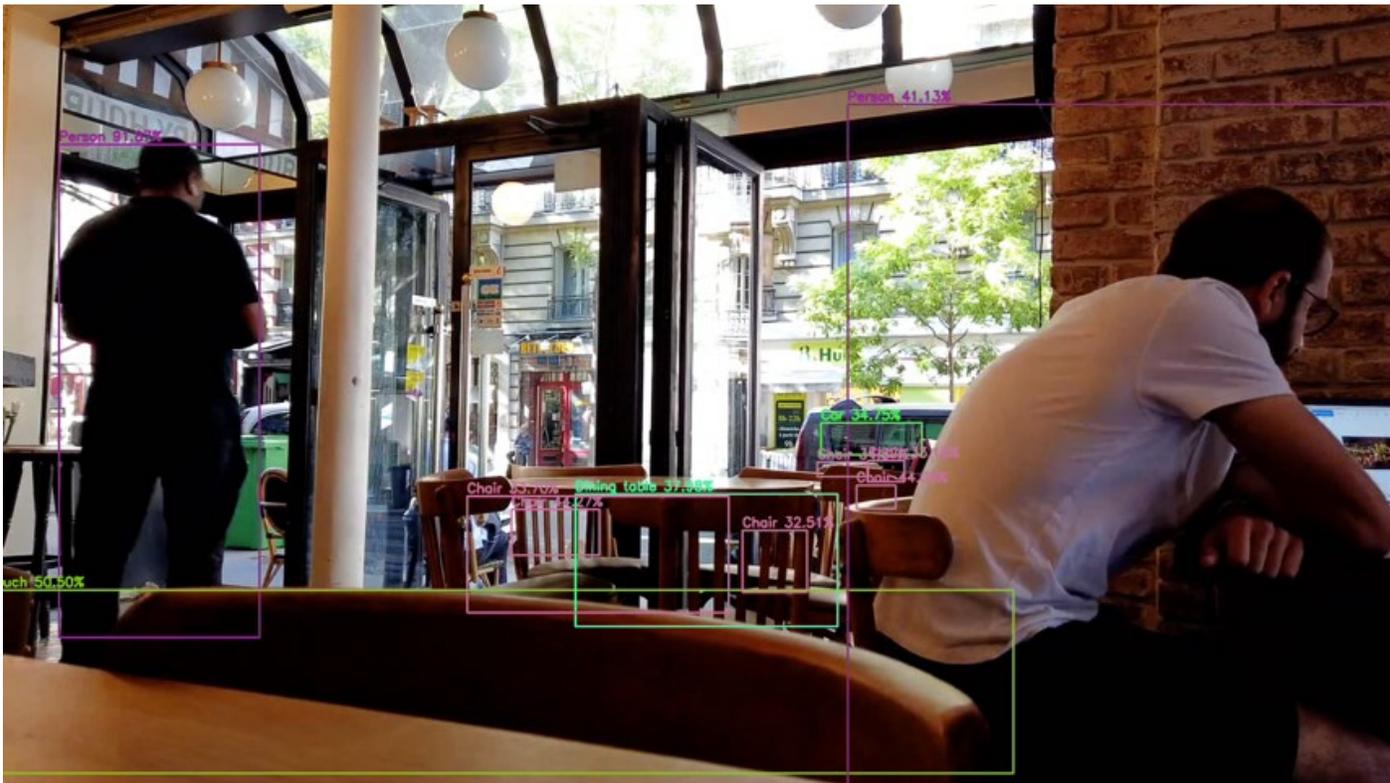
R-CNN is a relatively simple and straightforward approach. Several popular object detection models belong to the R-CNN family, including Fast R-CNN and Mask R-CNN.

5 PyArmNN object recognition

This section of the guide describes how to use PyArmNN to perform object recognition, using the [PyArmNN Object Detection Sample Application](#) as an example.

The sample application takes a model and video file or camera feed as input. The application then runs inference on each frame. Finally, the application draws bounding boxes around detected objects, with the corresponding labels and confidence scores overlaid.

The following image shows an example of one frame from the resulting video file with bounding boxes and confidence scores:



The PyArmNN Object Detection Sample Application performs the following steps:

1. Initialization
 - a. Read from the video source
 - b. Prepare labels and model-specific functions
2. Create a network
 - a. Create the parser and import a graph
 - b. Optimize the graph for the compute device
 - c. Create input and output binding information
3. Object detection pipeline

- a. Preprocess the captured frame
 - b. Make the input and output tensors
 - c. Execute inference
4. Postprocessing
 - a. Decode and process the inference output
 - b. Draw the bounding boxes
 - c. Run the application

The following subsections describe these steps.

5.1 Read from the video source

The application parses the supplied user arguments and loads the specified video file or stream into an OpenCV `cv2.VideoCapture` object. The application then uses this object to capture frames from the source with the `read()` function.

The `VideoCapture` object also provides information about the source, like the frame rate and resolution of the input video. The application uses this information to create a `cv2.VideoWriter` object. This object is used at the end of every loop to write the processed frame to an output video file.

5.2 Prepare labels and model-specific functions

To interpret the inference result on the loaded network, an application must load the labels that are associated with the model. In the sample application, the `dict_labels()` function creates a dictionary that is keyed on the classification index at the output node of the model. The values in the dictionary map each label to a randomly generated RGB color. This mapping means that each class has a unique color, which is useful when plotting the bounding boxes of detected objects in a frame.

The user-specified model accesses and returns functions to decode and process the inference output, along with a resize factor. This resize factor is used when plotting bounding boxes, to ensure that they are scaled to their correct position in the original frame.

5.3 Create the parser and import a graph

A PyArmNN application must import a graph from file using an appropriate parser. Arm NN provides parsers for various model file types, including TFLite, TF, and ONNX. These parsers are libraries for loading neural networks of various formats into the Arm NN runtime.

Because both the Yolo v3 and SSD models are in the TFLite format, the sample application uses the TFLite parser `armnnTfLiteParser` to process the models.

The `CreateNetworkFromBinaryFile()` function creates the parser and loads the network file. The parser then constructs the underlying Arm NN graph from the network file.

5.4 Optimize the graph for the compute device

Arm NN supports optimized execution on multiple CPU, GPU, and Ethos-N NPU devices. Before executing a graph, the application must select the appropriate device context by using `IRuntime()` to create a runtime context with default options.

We can optimize the imported graph by specifying a list of backends in order of preference and implementing backend-specific optimizations. A unique string identifies each one of these backends. For example:

- `CpuAcc` represents the CPU backend.
- `GpuAcc` represents the GPU backend.
- `CpuRef` represents the CPU reference kernels.

Arm NN splits the entire graph into subgraphs based on these backends. Each subgraph is then optimized, and the corresponding subgraph in the original graph is substituted with its optimized version.

The `Optimize()` function optimizes the graph for inference, then `LoadNetwork()` loads the optimized network onto the compute device. The `LoadNetwork()` function also creates the backend-specific workloads for the layers and a backend-specific workload factory.

5.5 Create input and output binding information

Parsers extract the input information for the network. The `GetSubgraphInputTensorNames()` function extracts all the input names and the `GetNetworkInputBindingInfo()` function obtains the input binding information of the graph.

The input binding information contains all the essential information about the input. This information is a tuple consisting of:

- Integer identifiers for bindable layers
- Tensor information including:
 - Data type
 - Quantization information
 - Number of dimensions
 - Total number of elements

Similarly, we can get the output binding information for an output layer by using the parser to retrieve output tensor names and calling the `GetNetworkOutputBindingInfo()` function.

5.6 Preprocess the captured frame

Each frame that is captured from the video source is read as a `ndarray` in BGR format. Each frame must then be preprocessed before being passed into the network.

This preprocessing step consists of the following:

1. Swap channels. In this example, swap BGR to RGB.
2. Resize the frame to the required resolution
3. Expand the dimensions of the array and perform data type conversion to match the model input layer.

You can read `input_binding_info` to obtain information about the shape and the data type of the input tensor. For example, SSD MobileNet V1 takes an input tensor with shape `[1, 300, 300, 3]` and data type `uint8`.

5.7 Make the input and output tensors

The `make_input_tensors()` function produces the input workload tensors.

The `make_output_tensors()` function produces the output workload tensors.

5.8 Execute inference

After creating the workload tensors, the compute device performs inference for the loaded network using the `EnqueueWorkload()` function of the runtime context. Calling the `workload_tensors_to_ndarray()` function obtains the inference results as a list of `ndarrays`.

5.9 Decode and process the inference output

The output from inference must be decoded to obtain information about detected objects in the frame.

The examples includes implementations of two networks, but you can implement your own network decoding solution. For more information, see [Implementing Your Own Network](#).

For SSD MobileNet V1 models, the application decodes the results to obtain the bounding box positions, classification index, confidence, and number of detections in the input frame.

For YOLO v3 tiny models, the application decodes the output and performs non-maximum suppression. This suppression filters out any weak detections below a confidence threshold and any redundant bounding boxes above an intersection-over-union (IoU) threshold.

Experiment with different threshold values for confidence and IoU to achieve the best visual results.

Detection results are returned as a list with the following form:

```
[class index, [box positions], confidence score]
```

Where [box positions] contains bounding box coordinates in the following form:

```
[x_min, y_min, x_max, y_max]
```

5.10 Draw the bounding boxes

The `draw_bounding_boxes()` function takes the inference results and draws bounding boxes around detected objects. This function also adds the associated label and confidence score. The labels dictionary that we created in [Preparing labels and model-specific functions](#) uses the class index of the detected object as a key to return the associated label and color for that class. The resize factor that we defined in [Preparing labels and model-specific functions](#) scales the bounding box coordinates to their correct positions in the original frame.

The processed frames are then written to file or displayed in a separate window.

5.11 Run the application

To run the video file with the Yolo v3 model with PyArmNN, use the following command:

```
python3 run_video_file.py -video_file_path <your_video> --model_file_path  
yolo_v3_tiny_darknet_fp32.tflite --model_name yolo_v3_tiny
```

To run the SSD model with PyArmNN, use the following command:

```
python3 run_video_file.py --video_file_path <your_video> --model_file_path  
ssd_mobilenet_v1.tflite --model_name ssd
```

6 Arm NN C++ API object recognition

This section of the guide describes how to use the Arm NN public C++ API to perform object recognition. We use the [Arm NN Object Detection Example](#) to illustrate the process.

The sample application takes a model and video file or camera feed as input. The application then runs inference on each frame. Finally, the application draws bounding boxes around detected objects, with the corresponding labels and confidence scores overlaid.

The Arm NN Object Detection Example performs the following steps:

1. Initialization
 - a. Read from the video source
 - b. Prepare labels and model-specific functions
2. Create a network
 - a. Create the parser and import a graph
 - b. Optimize the graph for the compute device
 - c. Create input and output binding information
3. Object detection pipeline
 - a. Preprocess the captured frame
 - b. Make the input and output tensors
 - c. Execute inference
4. Postprocessing
 - a. Decode and process the inference output
 - b. Draw the bounding boxes
 - c. Run the application

The following subsections describe these steps.

6.1 Read from the video source

After parsing user arguments, the application loads the chosen video file or stream into an OpenCV `cv::VideoCapture` object. The main function uses the [IFrameReader](#) interface and the OpenCV-specific implementation [CvVideoFrameReader](#) to capture frames from the source using the `ReadFrame()` function.

The `CvVideoFrameReader` object also provides information about the input video. Based on this information and the application arguments, the application creates one of the implementations of the [IFrameOutput](#) interface: `CvVideoFileWriter` or `CvWindowOutput`. The created object is used at the end of every loop to do one of the following:

- `CvVideoFileWriter` uses `cv::VideoWriter` with an `ffmpeg` backend to write the processed frame to an output video file.
- `CvWindowOutput` uses the `cv::imshow()` function to write the processed frame to a GUI window. See the `GetFrameSourceAndSink` function in `Main.cpp` for more details.

6.2 Prepare labels and model-specific functions

To interpret the result of running inference on the loaded network, the application must load the labels associated with the model. In the provided example code, the `AssignColourToLabel` function creates a vector of `[label, color]` pairs. The vector is ordered according to the object class index at the output node of the model. Labels are assigned with a randomly generated RGB color. This ensures that each class has a unique color which is helpful when plotting the bounding boxes of various detected objects in a frame.

Depending on the model that is being used, the `CreatePipeline` function returns a specific implementation of the object detection pipeline.

6.3 Create a network

All operations with Arm NN and networks are encapsulated in `ArmnnNetworkExecutor` class.

6.4 Create the parser and import a graph

Arm NN SDK imports the graph from a file using the appropriate parser.

Arm NN SDK provides parsers for reading graphs from various model formats. The example application focuses on the `.tflite`, `.pb`, and `.onnx` models.

Based on the extension of the provided model file, the corresponding parser is created, and the network file loaded with the `CreateNetworkFromBinaryFile()` method. The parser creates the underlying Arm NN graph.

The example application accepts `.tflite` format model files, using `ITfLiteParser`:

```
#include "armnnTfLiteParser/ITfLiteParser.hpp"

armnnTfLiteParser::ITfLiteParserPtr parser =
    armnnTfLiteParser::ITfLiteParser::Create();

armnn::INetworkPtr network = parser->CreateNetworkFromBinaryFile(modelPath.c_str());
```

6.5 Optimize the graph for the compute device

Arm NN supports optimized execution on multiple CPU and GPU devices. Before executing a graph, the application must select the appropriate device context. The example application creates a runtime context with default options with `IRuntime()`, as shown in the following code:

```
#include "armnn/ArmNN.hpp"
```

```
auto runtime = armnn::IRuntime::Create(armnn::IRuntime::CreationOptions());
```

The application optimizes the imported graph by specifying a list of backends in order of preference and implementing backend-specific optimizations. A unique string identifies each of the backends, for example `CpuAcc`, `GpuAcc`, `CpuRef`.

For example, the example application specifies backend optimizations, as shown in the following code:

```
std::vector<armnn::BackendId> backends{"CpuAcc", "GpuAcc", "CpuRef"};
```

Internally and transparently, Arm NN splits the graph into subgraphs based on the specified backends. Arm NN optimizes each of the subgraphs and, if possible, substitutes the corresponding subgraph in the original graph with its optimized version.

The application uses the `Optimize()` function to optimize the graph for inference, then loads the optimized network onto the compute device with `LoadNetwork()`. The `LoadNetwork()` function creates:

- The backend-specific workloads for the layers
- A backend-specific workload factory which creates the workloads.

The example application contains the following code:

```
armnn::IOptimizedNetworkPtr optNet = Optimize(*network,
                                             backends,
                                             m_Runtime->GetDeviceSpec(),
                                             armnn::OptimizerOptions());

std::string errorMessage;

runtime->LoadNetwork(0, std::move(optNet), errorMessage);

std::cerr << errorMessage << std::endl;
```

6.6 Create input and output binding information

Parsers can also extract input information for the network. The application calls `GetSubgraphInputTensorNames` to extract all the input names, then `GetNetworkInputBindingInfo` binds the input points of the graph. The example application contains the following code:

```
std::vector<std::string> inputNames = parser->GetSubgraphInputTensorNames(0);
auto inputBindingInfo = parser->GetNetworkInputBindingInfo(0, inputNames[0]);
```

The input binding information contains all the essential information about the input.

This information is a tuple consisting of:

- Integer identifiers for bindable layers
- Tensor information including:
 - Data type
 - Quantization information
 - Number of dimensions
 - Total number of elements

Similarly, the application gets the output binding information for an output layer by using the parser to retrieve output tensor names and calling `GetNetworkOutputBindingInfo()`.

6.7 Object detection pipeline

The generic object detection pipeline contains the following three steps:

1. Perform data pre-processing.
2. Run inference.
3. Decode inference results in the post-processing step.

See [ObjDetectionPipeline](#) and the implementations for [MobileNetSSDv1](#) and [YoloV3Tiny](#) for more details.

6.8 Preprocess the captured frame

The application reads each frame captured from source as a `cv::Mat` in BGR format. The channels are swapped to RGB in frame reader code, as follows:

```
cv::Mat processed;
...
objectDetectionPipeline->PreProcessing(frame, processed);
```

The preprocessing step consists of resizing the frame to the required resolution, padding, and converting the data types to match the model input layer. For example, the example application uses SSD MobileNet V1 which takes an input tensor with shape `[1, 300, 300, 3]` and data type `uint8`.

The preprocessing step returns a `cv::Mat` object containing data ready for inference.

6.9 Execute inference

The following code shows how the application executes inference:

```
od::InferenceResults results;
...
objectDetectionPipeline->Inference(processed, results);
```

The inference step calls the `ArmnnNetworkExecutor::Run` method that prepares input tensors and executes inference. A compute device performs inference for the loaded network using the `EnqueueWorkload()` function of the runtime context. For example:

```
//const void* inputData = ...;
//outputTensors were pre-allocated before

armnn::InputTensors inputTensors = {{
    inputBindingInfo.first, armnn::ConstTensor(inputBindingInfo.second, inputData)}};
runtime->EnqueueWorkload(0, inputTensors, outputTensors);
```

The application allocates memory for output data once and maps it to output tensor objects. After successful inference, the application reads data from the pre-allocated output data buffer. See [ArmnnNetworkExecutor::ArmnnNetworkExecutor](#) and [ArmnnNetworkExecutor::Run](#) for more information.

6.10 Decode and process the inference output

The application must decode the output from inference to obtain information about the detected objects in the frame. The example application contains implementations for two networks, or you can implement your own network decoding solution.

For SSD MobileNet V1 models, the application decodes the results to obtain the bounding box positions, classification index, confidence, and number of detections in the input frame. See [SSDResultDecoder](#) for more details.

For YOLO V3 Tiny models, the application decodes the output and performs non-maximum suppression. This suppression filters out weak detections below a confidence threshold and any redundant bounding boxes above an intersection-over-union (IoU) threshold. See [YoloResultDecoder](#) for more details.

Experiment with different threshold values for confidence and IoU to achieve the best visual results.

The detection results are always returned as a vector of `DetectedObject`, with the box positions list containing bounding box coordinates in the following form:

```
[x_min, y_min, x_max, y_max]
```

6.11 Draw the bounding boxes

The post-processing step accepts a callback function which is invoked when decoding finishes. The application uses this callback function to draw detections on the initial frame. The example application uses the output detections and the `AddInferenceOutputToFrame` function to draw bounding boxes around detected objects and add the associated label and confidence score. The following code shows the post-processing step in detail:

```
objectDetectionPipeline->PostProcessing(results,  
    [&frame, &labels](od::DetectedObjects detects) -> void {  
        AddInferenceOutputToFrame(detects, *frame, labels);  
    });
```

The processed frames are written to a file or displayed in a separate window.

6.12 Run the application

After building the application executable, you can run it with the following command-line options:

- `--video-file-path`: Specifies the path to the video file. This option is required.
- `--model-file-path`: Specifies the path to the object detection model. This option is required.
- `--label-path`: Specifies the path to the label set for the model file. This option is required.

- `--model-name`: Specifies the name of the model used for object detection. Valid values are `SSD_MOBILE` and `YOLO_V3_TINY`. This option is required.
- `--output-video-file-path`: Specifies the path to the output video file. This is optional. The default is `/tmp/output.avi`.
- `--preferred-backends`: Specifies the backends in preference order, separated by a comma. Valid values include `CpuAcc`, `CpuRef`, and `GpuAcc`. This is optional. The default is `CpuRef`, the reference kernel on CPU.
- `--help`: Displays all the available command-line options.

To run object detection on a video file and output the result to another video file, use the following commands:

```
LD_LIBRARY_PATH=/path/to/armnn/libs:/path/to/opencv/libs
./object_detection_example --label-path /path/to/labels/file
--video-file-path /path/to/video/file --model-file-path /path/to/model/file
--model-name [YOLO_V3_TINY | SSD_MOBILE]
--output-video-file-path /path/to/output/file
```

To run object detection on a video file and output the result to a window GUI, use the following commands:

```
LD_LIBRARY_PATH=/path/to/armnn/libs:/path/to/opencv/libs
./object_detection_example --label-path /path/to/labels/file
--video-file-path /path/to/video/file --model-file-path /path/to/model/file
--model-name [YOLO_V3_TINY | SSD_MOBILE]
```

7 Related information

Here are some resources related to material in this guide:

- Joseph Redmon, Ali Farhadi. [YOLOv3: An incremental improvement](#)
- Joseph Redmon, [You only look once: Unified, real-time object detection](#). CVPR 2016
- Wei Liu et al. [SSD: Single Shot MultiBox Detector](#). ECCV 2016

8 Next steps

In this guide, we examined sample applications that use the Arm NN C++ API and the PyArmNN Python extension to create a real-time object detection system.

You can now use the knowledge that you have gained to write your own applications.

To learn more, you can read our many [guides related to AI and Machine Learning](#) to learn more about PyArmNN and Arm NN.